# Duality in RKHSs with Infinite Dimensional Outputs: Application to Robust Losses

Pierre Laforgue [1]    Alex Lambert [1]    Luc Brogat-Motte [1]    Florence d'Alché-Buc [1]

## Abstract

Operator-Valued Kernels (OVKs) and associated vector-valued Reproducing Kernel Hilbert Spaces provide an elegant way to extend scalar kernel methods when the output space is a Hilbert space. Although primarily used in finite dimension for problems like multi-task regression, the ability of this framework to deal with infinite dimensional output spaces unlocks many more applications, such as functional regression, structured output prediction, and structured data representation. However, these sophisticated schemes crucially rely on the kernel trick in the output space, so that most of previous works have focused on the square norm loss function, completely neglecting robustness issues that may arise in such surrogate problems. To overcome this limitation, this paper develops a duality approach that allows to solve OVK machines for a wide range of loss functions. The infinite dimensional Lagrange multipliers are handled through a *Double Representer Theorem*, and algorithms for $\epsilon$-insensitive losses and the Huber loss are thoroughly detailed. Robustness benefits are emphasized by a theoretical stability analysis, as well as empirical improvements on structured data applications.

## 1. Introduction

Due to increasingly available streaming and network data, learning to predict complex objects such as structured outputs or time series has attracted a great deal of attention in machine learning. Extending the well known kernel methods devoted to non-vectorial data (Hofmann et al., 2008), several kernel-based approaches have emerged to deal with complex output data. While Structural SVM

and variants cope with discrete structures (Tsochantaridis et al., 2005; Joachims et al., 2009) through structured losses, Operator-Valued Kernels (OVKs) and vector-valued Reproducing Kernel Hilbert Spaces (vv-RKHSs, Micchelli and Pontil (2005); Carmeli et al. (2006; 2010)) provide a unique framework to handle both functional and structured outputs. Vv-RKHSs are classes of functions that map an arbitrary input set $\mathcal{X}$ to some output Hilbert space $\mathcal{Y}$ (Senkene and Tempel'man, 1973; Caponnetto et al., 2008). Primarily used with finite dimensional outputs ($\mathcal{Y} = \mathbb{R}^p$) to solve multi-task regression (Micchelli and Pontil, 2005; Baldassarre et al., 2012) and multiple class classification (Dinuzzo et al., 2011), OVK methods have further been exploited to handle outputs in infinite dimensional Hilbert spaces. This has unlocked numerous applications, such as functional regression (Kadri et al., 2010; 2016), structured prediction (Brouard et al., 2011; Kadri et al., 2013), infinite quantile regression (Brault et al., 2019), or structured data representation learning (Laforgue et al., 2019). Nonetheless, these sophisticated schemes often come along with a basic loss function: the output space squared norm, neglecting desirable properties such as parsimony and robustness.

In nonparametric modeling, model parsimony boils down to data sparsity, *e.g.* reducing the number of training data points on which the model relies to make a prediction. Such a property is highly valuable (Hastie et al., 2015): not only does it prevent overfitting but it also alleviates the inherent computational load of optimization and prediction, allowing to scale to larger datasets. Another appealing property of a regression tool is robustness to outliers (Huber, 1964; Zhu et al., 2008). Real data may suffer from incorrect feature measurements and spurious annotations, leading to training datasets contaminated with outliers. Then, minimizing the squared loss is inappropriate as the least-squares estimates behave poorly when the residuals distribution is not normal, but rather heavy-tailed. In (scalar) kernel methods, these two properties – data sparsity and robustness to outliers – are imposed through the choice of appropriate losses. Data sparsity is leveraged by using $\epsilon$-insensitive losses, exploited in the well known Support Vector Regression (Drucker et al., 1997) while robust regression (Fung and Mangasarian, 2000) can be obtained by minimizing the Huber loss function (Huber, 1964). Driven by three emblematic learning tasks,

structured prediction, functional regression, and structured data representation, we propose a general duality framework that enables sparse data regression and robust regression, even when working in vv-RKHSs with infinite-dimensional outputs. Although extensively used within scalar kernel methods, very few attempts have been made to adapt duality to vv-RKHSs. In Brouard et al. (2016b), dualization is presented, but only used in the maximum margin regression scenario. Sangnier et al. (2017) consider a wider class of loss functions, including $\epsilon$-insensitive losses to leverage data sparsity, but only in the case of matrix-valued kernels (Álvarez et al., 2012), for which the dual problem is finite dimensional. For a general OVK however, the dual problem is to be solved over $\mathcal{Y}^n$, and is intractable without additional work when $\mathcal{Y}$ is infinite dimensional. We first notice that the extensions of $\epsilon$-insensitive losses and the Huber loss to general Hilbert space are (still) expressed as convolutions of simpler losses whose Fenchel-Legendre (FL) transforms are known. Inspired by this remark, we identify general conditions on the OVKs and FL transforms to establish a *Double Representer Theorem* allowing to work with matrix parameterized representations. In particular, a careful use of the duality principle considerably broadens the range of loss functions for which OVK solutions are computable. The present work thus aims at developing a comprehensive methodology to solve these dual problems.

The rest of the paper is organized as follows. In Section 2, we introduce OVKs, recall the general formulation of dual problems for OVK machines, and derive their solvable finite dimensional reformulation. Section 3 is devoted to specific instantiations of this problem for $\epsilon$-insensitive losses and the Huber loss, with algorithms duly explicited. In Section 4, we apply our framework to induce sparsity and robustness into structured prediction, functional regression, and structured data representation. Proofs are postponed to the Appendix.

## 2. Learning in vv-RKHSs

After reminders on OVKs and vv-RKHS learning theory, this section exposes the duality approach for the regularized empirical risk minimization problem in vv-RKHSs. Two strategies are then detailed to solve infinite dimensional dual problems, either under an assumption on the kernel, or by approximating the dual. In the following, $\mathcal{Y}$ is assumed to be a separable Hilbert space.

**Definition 1.** *An OVK is an application $\mathcal{K}\colon \mathcal{X}\times\mathcal{X} \to \mathcal{L}(\mathcal{Y})$, that satisfies the following two properties for all $n \in \mathbb{N}^*$:*

1) $\forall (x,x') \in \mathcal{X}\times\mathcal{X}, \qquad \mathcal{K}(x,x') = \mathcal{K}(x',x)^\#,$

2) $\forall\, (x_i,y_i)_{i=1}^n \in (\mathcal{X}\times\mathcal{Y})^n,\ \sum_{i,j=1}^n \langle y_i, \mathcal{K}(x_i,x_j)y_j\rangle_{\mathcal{Y}} \geqslant 0,$

*with $\mathcal{L}(E)$ the set of bounded linear operators on vector space $E$, and $A^\#$ the adjoint of any operator $A$.*

A simple example of OVK is the *separable kernel*.

**Definition 2.** *$\mathcal{K}:\mathcal{X}\times\mathcal{X} \to \mathcal{L}(\mathcal{Y})$ is a* separable kernel *iff there exist a scalar kernel $k : \mathcal{X}\times\mathcal{X} \to \mathbb{R}$ and a positive semi-definite operator $A \in \mathcal{L}(\mathcal{Y})$ such that for all $(x,x') \in \mathcal{X}^2$ it holds: $\mathcal{K}(x,x') = k(x,x')A.$*

Similarly to scalar-valued kernels, an OVK can be uniquely associated to a functional space from $\mathcal{X}$ to $\mathcal{Y}$: its vv-RKHS.

**Theorem 1.** *Let $\mathcal{K}$ be an OVK, and for $x \in \mathcal{X}$, let $\mathcal{K}_x\colon y \mapsto \mathcal{K}_x y \in \mathcal{F}(\mathcal{X},\mathcal{Y})$ the linear operator such that: $\forall x' \in \mathcal{X},\ (\mathcal{K}_x y)(x') = \mathcal{K}(x',x)y$. Then, there is a unique Hilbert space $\mathcal{H}_{\mathcal{K}} \subset \mathcal{F}(\mathcal{X},\mathcal{Y})$ the vv-RKHS associated to $\mathcal{K}$ such that $\forall x \in \mathcal{X}$ it holds:*

*(i) $\mathcal{K}_x$ spans the space $\mathcal{H}_{\mathcal{K}}$ ($\forall y \in \mathcal{Y}\colon \mathcal{K}_x y \in \mathcal{H}_{\mathcal{K}}$)*

*(ii) $\mathcal{K}_x$ is bounded for the uniform norm*

*(iii) $\forall f \in \mathcal{H}_{\mathcal{K}},\ f(x) = \mathcal{K}_x^\# f$ (reproducing property)*

Given a sample $\mathcal{S} = \{(x_i,y_i)_{i=1}^n\} \in (\mathcal{X}\times\mathcal{Y})^n$ of $n$ i.i.d. realizations of a generic random variable $(X,Y)$, an OVK $\mathcal{K}: \mathcal{X}\times\mathcal{X} \to \mathcal{L}(\mathcal{Y})$, a convex loss function $\ell : \mathcal{Y}\times\mathcal{Y} \to \mathbb{R}$, and a regularization parameter $\Lambda > 0$, the general form of an OVK-based learning problem is to find $\hat{h}$ that solves:

$$\min_{h \in \mathcal{H}_{\mathcal{K}}}\ \frac{1}{n}\sum_{i=1}^n \ell(h(x_i), y_i) + \frac{\Lambda}{2}\|h\|_{\mathcal{H}_{\mathcal{K}}}^2. \qquad (1)$$

Similarly to scalar ones, a crucial tool in operator-valued kernel methods is the *Representer Theorem*, ensuring that $\hat{h}$ actually pertains to a reduced subspace of $\mathcal{H}_{\mathcal{K}}$.

**Theorem 2.** *(Theorem 4.2 in Micchelli and Pontil (2005)) There exists $(\hat{\alpha}_i)_{i=1}^n \in \mathcal{Y}^n$ such that*

$$\hat{h} = \frac{1}{\Lambda n}\sum_{i=1}^n \mathcal{K}(\cdot, x_i)\hat{\alpha}_i.$$

Although Theorem 2 drastically downscales the search domain (from $\mathcal{H}_{\mathcal{K}}$ to $\mathcal{Y}^n$), it gives no further information about the $(\hat{\alpha}_i)_{i=1}^n$. One way to gain insight about these coefficients is to perform Problem (1)'s dualization, with the notation $\ell_i : y \in \mathcal{Y} \mapsto \ell(y, y_i)$ for any $i \leqslant n$.

**Theorem 3.** *(Appendix B in Brouard et al. (2016b)) The solution to Problem (1) is given by*

$$\hat{h} = \frac{1}{\Lambda n}\sum_{i=1}^n \mathcal{K}(\cdot, x_i)\hat{\alpha}_i, \qquad (2)$$

*with $(\hat{\alpha}_i)_{i=1}^n \in \mathcal{Y}^n$ the solutions to the dual problem*

$$\min_{(\alpha_i)_{i=1}^n \in \mathcal{Y}^n}\ \sum_{i=1}^n \ell_i^\star(-\alpha_i) + \frac{1}{2\Lambda n}\sum_{i,j=1}^n \langle \alpha_i, \mathcal{K}(x_i,x_j)\alpha_j\rangle_{\mathcal{Y}}, \qquad (3)$$

*where $f^\star : \alpha \in \mathcal{Y} \mapsto \sup_{y \in \mathcal{Y}}\langle \alpha, y\rangle_{\mathcal{Y}} - f(y)$ denotes the Fenchel-Legendre transform of a function $f : \mathcal{Y} \to \mathbb{R}$.*

Refer to Appendix A.1 for Theorem 3's proof, that has been reproduced for self-containedness. Dualization brings in additional information about the optimal coefficients (notice nonetheless that Theorem 2 holds true for a much wider class of problems). As it is, Problem (3) is however of little interest, since the optimization must be performed on the infinite dimensional space $\mathcal{Y}^n$. Depending on the problem, we propose two solutions: either using a *Double Representer Theorem*, or by approximating Problem (3).

**Notation.** If $\mathcal{K}$ is identity decomposable (i.e. $\mathcal{K} = k \, \mathbf{I}_\mathcal{Y}$), $K^X$ and $K^Y$ denote the input and output gram matrices. For any matrix $M$, $M_{i:}$ represents its $i^{th}$ line, and $\|M\|_{p,q}$ its $\ell_{p,q}$ row wise mixed norm, *i.e.* the $\ell_q$ norm of the $\ell_p$ norms of its lines. $\chi_S$ denotes the characteristic function of a set $S$, null on $S$ and equal to $+\infty$ otherwise, $f \square g$ is the infimal convolution of $f$ and $g$ (Bauschke et al., 2011), $(f \square g)(x) = \inf_y f(y) + g(x - y)$. Finally, $\#S$ is the cardinality of any set $S$, and $\|\cdot\|_{\text{op}}$ the operator norm.

### 2.1. The Double Representer Theorem

In order to make Problem (3) solvable, we need assumptions on the loss and the kernel. Let $\mathbf{Y}$ denote $\text{span}(y_i, \ i \leqslant n)$. Assumptions 1 and 2 characterize admissible losses through conditions on their Fenchel-Legendre (FL) transforms. They are standard for kernel methods, and ensure computability by stipulating that only dot products are involved.

**Assumption 1.** $\forall i \leqslant n, \ \forall (\alpha^\mathbf{Y}, \alpha^\perp) \in \mathbf{Y} \times \mathbf{Y}^\perp$, it holds $\ell_i^\star(\alpha^\mathbf{Y}) \leqslant \ell_i^\star(\alpha^\mathbf{Y} + \alpha^\perp)$.

**Assumption 2.** $\forall i \leqslant n, \exists L_i : \mathbb{R}^{n+n^2} \to \mathbb{R}$ such that for all $\boldsymbol{\omega} = (\omega_j)_{j \leqslant n} \in \mathbb{R}^n$, $\quad \ell_i^\star(-\sum_{j=1}^n \omega_j \, y_j) = L_i(\boldsymbol{\omega}, K^Y)$.

Regarding the OVK, the key point is Assumption 3. Roughly speaking, $\mathbf{Y}$ is what we *see* and *know* about output space $\mathcal{Y}$, while $\mathbf{Y}^\perp$ represents the part we *ignore*. What we need is an OVK somewhat *aligned* with the outputs, in the sense that the little we know about $\mathcal{Y}$ should be preserved through $\mathcal{K}$. As for Assumption 4, it helps simplifying the computations.

**Assumption 3.** $\forall i, j \leqslant n$, $\mathbf{Y}$ is invariant by $\mathcal{K}(x_i, x_j)$, i.e. $\forall y \in \mathcal{Y}, \ y \in \mathbf{Y} \Rightarrow \mathcal{K}(x_i, x_j)y \in \mathbf{Y}$.

**Remark 1.** *It is important to notice that we do not need Assumption 3 to hold true for every collection $\{y_i\}_{i \leqslant n} \in \mathcal{Y}^n$. It rather constitutes an a posteriori condition to ensure that the kernel is aligned with the training sample at hand. If $\mathcal{Y}$ is finite dimensional, one may hope that with sufficiently many outputs, then $\mathbf{Y}$ spans $\mathcal{Y}$, and every matrix-valued kernel then fits. If $\mathcal{Y}$ is infinite dimensional, identity-decomposable kernels are admissible (which despite simple expressions may describe nontrivial dependences in infinite dimensional spaces). Moreover, separable kernels with operators similar to the empirical covariance $\sum_i y_i \otimes y_i$ (Kadri et al., 2013) are also eligible, opening the door to ad-hoc and learned kernels, see Appendix A.8 for further examples.*

**Assumption 4.** *There exist $T \geqslant 1$, and for every $t \leqslant T$ admissible scalar kernels $k_t : \mathcal{X} \times \mathcal{X} \to \mathbb{R}$ as well as positive semi-definite operators $A_t \in \mathcal{L}(\mathcal{Y})$, such that for all $(x, x') \in \mathcal{X}^2$ it holds: $\mathcal{K}(x, x') = \sum_{t=1}^T k_t(x, x') A_t$.*

Under Assumption 4, $K_t^X$ and $K_t^Y$ denote the matrices such that $[K_t^X]_{ij} = k_t(x_i, x_j)$, $[K_t^Y]_{ij} = \langle y_i, A_t y_j \rangle_\mathcal{Y}$. Notice that it is by no means restrictive, since every shift-invariant OVK can be approximated arbitrarily closely by kernels satisfying Assumption 4. Furthermore, if for all $t \leqslant T$, $A_t$ keeps $\mathbf{Y}$ invariant, then Assumption 3 is directly fulfilled. Under these assumptions, Theorem 4 proves that the optimal coefficients lie in $\mathbf{Y}^n$, ensuring the solutions computability.

**Theorem 4.** *Let $\ell : \mathcal{Y} \times \mathcal{Y} \to \mathbb{R}$ be a loss function with Fenchel-Legendre transforms satisfying Assumptions 1 and 2, and $\mathcal{K}$ be an OVK verifying Assumption 3. Then, the solution to Problem (1) is given by*

$$\hat{h} = \frac{1}{\Lambda n} \sum_{i,j=1}^n \mathcal{K}(\cdot, x_i) \, \hat{\omega}_{ij} \, y_j, \tag{4}$$

*with $\hat{\Omega} = [\hat{\omega}_{ij}] \in \mathbb{R}^{n \times n}$ the solution to the dual problem*

$$\min_{\Omega \in \mathbb{R}^{n \times n}} \ \sum_{i=1}^n L_i\left(\Omega_{i:}, K^Y\right) + \frac{1}{2\Lambda n} \mathbf{Tr}\left(\tilde{M}^\top (\Omega \otimes \Omega)\right),$$

*with $M$ the $n^4$ tensor such that $M_{ijkl} = \langle y_k, \mathcal{K}(x_i, x_j) y_l \rangle_\mathcal{Y}$, and $\tilde{M}$ its rewriting as a $n^2 \times n^2$ block matrix. If kernel $\mathcal{K}$ further satisfies Assumption 4, then tensor $M$ simplifies to $M_{ijkl} = \sum_{t=1}^T [K_t^X]_{ij} [K_t^Y]_{kl}$, and the problem rewrites*

$$\min_{\Omega \in \mathbb{R}^{n \times n}} \ \sum_{i=1}^n L_i\left(\Omega_{i:}, K^Y\right) + \frac{1}{2\Lambda n} \sum_{t=1}^T \mathbf{Tr}\left(K_t^X \Omega K_t^Y \Omega^\top\right). \tag{5}$$

See Appendix A.2 for the proof. This theorem can be seen as a *Double Representer Theorem*, since both theorems share analogous proofs and consequences: a search domain reduction, respectively from $\mathcal{H}_\mathcal{K}$ to $\mathcal{Y}^n$, and $\mathcal{Y}^n$ to $\mathbb{R}^{n \times n}$.

**Remark 2.** *The Double Representer Theorem emphasizes that only the knowledge of the $n^4$ tensor $M$ is required to make OVK problems in infinite dimensional output spaces computable. Although it might seem prohibitive at first sight, one has to keep in mind that, like for scalar kernel methods, a first $n^2$ cost is needed to use (input) kernels with infinite dimensional feature maps, while the second $n^2$ cost allows for handling infinite dimensional outputs. In the case of a decomposable kernel, one has $M_{ijkl} = K_{ij}^X K_{kl}^Y$. One only needs two $n^2$ matrices, recovering the scalar complexity.*

We now present a non-exhaustive list of admissible losses (one may refer to Appendix A.3 for the proof).

**Proposition 1.** *The following losses have Fenchel-Legendre transforms verifying Assumptions 1 and 2:*

- $\ell_i(y) = f(\langle y, z_i \rangle)$, $z_i \in Y$ and $f : \mathbb{R} \to \mathbb{R}$ *convex. This encompasses maximum-margin regression, obtained with $z_i = y_i$ and $f(t) = \max(0, 1 - t)$.*

- $\ell(y) = f(\|y\|)$, $f : \mathbb{R}_+ \to \mathbb{R}$ *convex increasing s.t. $t \mapsto \frac{f'(t)}{t}$ is continuous over $\mathbb{R}_+$. This includes all power functions $\frac{\lambda}{\eta}\|y\|_{\mathcal{Y}}^\eta$ for $\eta > 1$ and $\lambda > 0$.*

- $\forall \lambda > 0$, *with $\mathcal{B}_\lambda$ the centered ball of radius $\lambda$,*

  - $\ell(y) = \lambda\|y\|$,    ▪ $\ell(y) = \lambda\|y\| \log(\|y\|)$,
  - $\ell(y) = \chi_{\mathcal{B}_\lambda}(y)$,    ▪ $\ell(y) = \lambda(\exp(\|y\|) - 1)$.

- $\ell_i(y) = f(y - y_i)$, $f^\star$ *verifying Assumptions 1 and 2.*

- *Any infimal convolution involving functions satisfying Assumptions 1 and 2. This encompasses $\epsilon$-insensitive losses (Sangnier et al., 2017), the Huber loss (Huber, 1964), and generally all Moreau or Pasch-Hausdorff envelopes (Moreau, 1962; Bauschke et al., 2011).*

### 2.2. Approximating the Dual Problem

If Assumption 3 is not satisfied, another way to get a finite dimensional decomposition similar to that of Theorem 4 is to approximate the dual problem. This may be done by restricting the dual variables to suitable finite dimensional subsets of $\mathcal{Y}$, if the following hypothesis on kernel $\mathcal{K}$ holds.

**Assumption 5.** *The kernel $\mathcal{K} = k \cdot A$ is a separable OVK, with $A$ a compact operator.*

Recalling that $A$ is by design self adjoint and positive, its compactness then allows for a spectral decomposition: there exists an orthonormal basis $(\psi_j)_{j=1}^\infty$ of $\mathcal{Y}$, and some positive $(\lambda_j)_{j=1}^\infty$, ordered in a non-increasing fashion and converging to zero, such that $A = \sum_{j=1}^\infty \lambda_j \psi_j \otimes \psi_j$ (Osborn, 1975).

Using such a basis, one can say that there exists $(\hat{\omega}_i)_{i=1}^n \in \ell^2(\mathbb{R})^n$ such that $\forall i \leqslant n, \hat{\alpha}_i = \sum_{j=1}^\infty \hat{\omega}_{ij}\psi_j$. Since this leads to an infinite size representation of the dual variables, the idea is then to restrict the search space to the eigenspace associated to the $m$ largest eigenvalues of $A$, for some $m > 0$. Let $\widetilde{\mathcal{Y}}_m$ denote span$(\{\psi_j\}_{j=1}^m)$, and $S = \text{diag}(\lambda_j)_{j=1}^m$. An approximated dual problem reads

$$\min_{(\alpha_i)_{i=1}^n \in \widetilde{\mathcal{Y}}_m^n} \sum_{i=1}^n \ell_i^\star(-\alpha_i) + \frac{1}{2\Lambda n} \sum_{i,j=1}^n \langle \alpha_i, \mathcal{K}(x_i, x_j)\alpha_j \rangle_{\mathcal{Y}}, \quad (6)$$

We now state a condition similar to Assumption 2, which makes the solution to Problem (6) computable.

**Assumption 6.** $\forall i \leqslant n, \exists L_i : \mathbb{R}^{2m} \to \mathbb{R}$ *such that $\forall \boldsymbol{\omega} = (\omega_j)_{j \leqslant m} \in \mathbb{R}^m$, $\ell_i^\star(-\sum_{j=1}^m \omega_j \psi_j) = L_i(\boldsymbol{\omega}, R_{i:})$, with $R \in \mathbb{R}^{n \times m}$ the matrix such that $R_{ij} = \langle y_i, \psi_j \rangle_{\mathcal{Y}}$.*

**Remark 3.** *Assumption 6 is similar to Assumption 2, except that the output Gram matrix $K^Y$ is replaced by matrix $R$ storing the dot products between the orthonormal family $\{\psi_j\}_{j=1}^m$ and the outputs. In particular, all losses explicited in Proposition 1 have FL transforms verifying Assumption 6.*

**Theorem 5.** *Let $\mathcal{K}$ be an OVK meeting Assumption 5 and $\ell : \mathcal{Y} \times \mathcal{Y} \to \mathbb{R}$ be a loss function with FL transforms satisfying Assumption 6. Then, Problem (6) is equivalent to*

$$\min_{\Omega \in \mathbb{R}^{n \times m}} \sum_{i=1}^n L_i(\Omega_{i:}, R_{i:}) + \frac{1}{2\Lambda n}\mathbf{Tr}\left(K^X \Omega S \Omega^\top\right). \quad (7)$$

*Denoting by $\hat{\Omega} = [\hat{\omega}_{ij}] \in \mathbb{R}^{n \times m}$ the solution to Problem (7), the associated predictor is finally given by*

$$\hat{h} = \frac{1}{\Lambda n} \sum_{i=1}^n \sum_{j=1}^m k(\cdot, x_i) \lambda_j \hat{\omega}_{ij} \psi_j, \quad (8)$$

**Remark 4.** *The rationale behind the above approximation is that under compactness of $A$, Equation (8) constitutes a reasonable approximation of Equation (2). Notice that Kadri et al. (2016) use a truncated spectral decomposition of the operator to implement a functional version of Kernel Ridge Regression, without resorting to dualization however.*

## 3. Application to Robust Losses

We now instantiate Theorem 4's dual problem for three loss functions encouraging data sparsity and robustness. They write as infimal convolutions, and are thus hardly tractable in the primal. Their dual problems enjoy simple resolution algorithms that are thoroughly detailed. A stability analysis is also carried out to highlight the hyperparameters impact.

### 3.1. Complete Dual Resolution for Three Robust Losses

As a first go, we recall the important notion of $\epsilon$-insensitive losses. Following in the footsteps of Sangnier et al. (2017), we extend them in a natural way from $\mathbb{R}^p$ to any Hilbert space $\mathcal{Y}$. To avoid additional notation, in this subsection $\ell$ denotes the loss taken w.r.t. one argument (previously $\ell_i$).

**Definition 3.** *Let $\ell : \mathcal{Y} \to \mathbb{R}_+$ be a convex loss such that $\ell(0) = 0$, and $\epsilon > 0$. The $\epsilon$-insensitive version of $\ell$, denoted $\ell_\epsilon$, is defined by $\ell_\epsilon(y) = (\ell \,\square\, \chi_{\mathcal{B}_\epsilon})(y)$, or again:*

$$\forall y \in \mathcal{Y}, \ \ell_\epsilon(y) = \begin{cases} 0 & \text{if } \|y\|_{\mathcal{Y}} \leqslant \epsilon \\ \inf_{\|d\|_{\mathcal{Y}} \leqslant 1} \ell(y - \epsilon d) & \text{otherwise} \end{cases}.$$

In other terms, $\ell_\epsilon(y)$ is the smallest value of $\ell$ within the ball of radius $\epsilon$ centered at $y$. As revealed by the next definition, natural choices for $\ell$ yield extensions of celebrated scalar loss functions to infinite dimensional Hilbert spaces.

**Definition 4.** *If $\ell = \| \cdot \|_{\mathcal{Y}}$, then $\| \cdot \|_{\mathcal{Y},\epsilon} = \max(\| \cdot \|_{\mathcal{Y}} - \epsilon, 0)$, and the related problem is the natural extension of $\epsilon$-SVR.*

*If $\ell = \| \cdot \|_{\mathcal{Y}}^2$, then $\| \cdot \|_{\mathcal{Y},\epsilon} = \max(\| \cdot \|_{\mathcal{Y}} - \epsilon, 0)^2$, and the related problem is called the $\epsilon$-insensitive Ridge regression.*

The third framework that nicely falls into our resolution methodology is the Huber loss regression ([Huber, 1964]). Tailored to induce robustness, the Huber loss function does not feature convolution with $\chi_{\mathcal{B}_\epsilon}$ but rather between the first two powers of the Hilbert norm (that used in Definition 4).

**Definition 5.** *The Huber loss of parameter $\kappa$ is given by $\ell_{H,\kappa}(y) = (\kappa \| \cdot \|_{\mathcal{Y}} \square \frac{1}{2} \| \cdot \|_{\mathcal{Y}}^2)(y)$, or again:*

$$\forall y \in \mathcal{Y}, \ \ell_{H,\kappa}(y) = \begin{cases} \frac{1}{2} \|y\|_{\mathcal{Y}}^2 & \text{if } \|y\|_{\mathcal{Y}} \leqslant \kappa \\ \kappa \left( \|y\|_{\mathcal{Y}} - \frac{\kappa}{2} \right) & \text{otherwise} \end{cases}.$$

Due to its asymptotic behavior as $\| \cdot \|_{\mathcal{Y}}$, the Huber loss is useful when the training data is heavy tailed or contains outliers. Illustrations of Definitions 4 and 5's loss functions in one and two dimensions are available in Appendix B. Interestingly, Problem (5) for these three losses – and an identity decomposable kernel – admits a very nice writing, allowing for an efficient resolution.

**Theorem 6.** *If $\mathcal{K} = k\,\mathbf{I}_{\mathcal{Y}}$, the solutions to the $\epsilon$-Ridge regression, $\kappa$-Huber regression, and $\epsilon$-SVR primal problems*

$$(P1) \quad \min_{h \in \mathcal{H}_{\mathcal{K}}} \ \frac{1}{2n} \sum_{i=1}^{n} \|h(x_i) - y_i\|_{\mathcal{Y},\epsilon}^2 + \frac{\Lambda}{2} \|h\|_{\mathcal{H}_{\mathcal{K}}}^2,$$

$$(P2) \quad \min_{h \in \mathcal{H}_{\mathcal{K}}} \ \frac{1}{n} \sum_{i=1}^{n} \ell_{H,\kappa}(h(x_i) - y_i) + \frac{\Lambda}{2} \|h\|_{\mathcal{H}_{\mathcal{K}}}^2,$$

$$(P3) \quad \min_{h \in \mathcal{H}_{\mathcal{K}}} \ \frac{1}{n} \sum_{i=1}^{n} \|h(x_i) - y_i\|_{\mathcal{Y},\epsilon} + \frac{\Lambda}{2} \|h\|_{\mathcal{H}_{\mathcal{K}}}^2,$$

*are given by Equation (4), with $\hat{\Omega} = \hat{W} V^{-1}$, and $\hat{W}$ the solution to the respective finite dimensional dual problems*

$$(D1) \quad \min_{W \in \mathbb{R}^{n \times n}} \ \frac{1}{2} \|AW - B\|_{\text{Fro}}^2 + \epsilon \|W\|_{2,1},$$

$$(D2) \quad \min_{W \in \mathbb{R}^{n \times n}} \ \frac{1}{2} \|AW - B\|_{\text{Fro}}^2,$$
$$\text{s.t.} \quad \|W\|_{2,\infty} \leqslant \kappa,$$

$$(D3) \quad \min_{W \in \mathbb{R}^{n \times n}} \ \frac{1}{2} \|AW - B\|_{\text{Fro}}^2 + \epsilon \|W\|_{2,1},$$
$$\text{s.t.} \quad \|W\|_{2,\infty} \leqslant 1,$$

*with $V$, $A$, $B$ such that: $VV^\top = K^Y$, $A^\top A = \frac{K^X}{\Lambda n} + \mathbf{I}_n$ (or $A^\top A = K^X/(\Lambda n)$ for the $\epsilon$-SVR), and $A^\top B = V$.*

Theorem 6's proof is detailed in Appendix A.5. If $\mathcal{K}$ is not identity decomposable, but only satisfies Assumption 4, the dual problems do not admit compact writings such as those of Theorem 6. Nonetheless, they are still easily solvable, and the standard Ridge regression is recovered for $\epsilon = 0$ or $\kappa = +\infty$. This is discussed at length in the Appendix.

Problem (D1) is a Multi-Task Lasso problem ([Obozinski et al., 2010]). It can be solved by Projected Gradient Descent (PGD), that involves the Block Soft Thresholding operator such that $\text{BST}(x, \tau) = (1 - \tau/\|x\|)_+ x$. Problem (D2) is a constrained least square problem, that also admits a resolution through PGD, but with the Projection operator such that $\text{Proj}(x, \tau) = \min(\tau/\|x\|, 1) x$. Finally, Problem (D3) combines both non-smooth terms and consequently both projection steps. Given a stepsize $\eta$, and $T$ a number of epoch, the algorithms are detailed in Algorithm 1. Note that $\tilde{K}$'s Singular Value Decomposition is not necessary, since the computations only involve $A^\top A = \tilde{K}$ and $A^\top B = V$.

---

**Algorithm 1** Projected Gradient Descents (PGDs)

---

**input** : Gram matrices $K^X$, $K^Y$, parameters $\Lambda$, $\epsilon$, $\kappa$

**init** : $\tilde{K} = \frac{1}{\Lambda n} K^X + \mathbf{I}_n$ (or $\tilde{K} = \frac{1}{\Lambda n} K^X$ for $\epsilon$-SVR),
$K^Y = VV^\top$, $W = \mathbf{0}_{\mathbb{R}^{n \times n}}$

**for** *epoch from* 1 *to* $T$ **do**

    // gradient step
    $W = W - \eta(\tilde{K}W - V)$
    // projection step
    **for** *row* $i$ *from* 1 *to* $n$ **do**
        $W_{i:} = \text{BST}(W_{i:}, \epsilon)$   // if Ridge or SVR
        $W_{i:} = \text{Proj}(W_{i:}, \kappa \text{ or } 1)$ // if Huber or SVR

**return** $W$

---

### 3.2. Approximate Dual Resolution with Huber Loss

In this section we solve Problem (6) for the Huber loss and $\mathcal{Y} = L^2[\Theta, \mu]$, with $\Theta$ a compact set endowed with measure $\mu$. A classical choice of OVK is then $\mathcal{K} = k_{\mathcal{X}} \cdot T_k$, $k_{\mathcal{X}}$ being a scalar kernel over the inputs, and $T_k$ the integral operator associated to a scalar kernel $k \colon \Theta \times \Theta \to \mathbb{R}$ defined for all $g \in L^2[\Theta, \mu]$ by $T_k g = \int_\Theta k(\cdot, \theta)g(\theta)\mathrm{d}\mu(\theta)$. Continuity of $k$ grants compactness of $T_k$, allowing for the methodology presented in Section 2.2. In the following, $(\lambda_j, \psi_j)_{j=1}^m$ denotes the eigendecomposition of $T_k$, which is dependent both in $k$ and $\mu$, and can be obtained by solving a differential equation derived from the eigenvalue problem. However, given that the optimal kernel $k$ is unknown, one can choose a Hilbertian basis $\{\psi_j\}_{j=1}^\infty$ of $L^2[\Theta, \mu]$ and a non-increasing summable sequence $(\lambda_j)_{j=1}^\infty \in \mathbb{R}_+^*$ to construct the kernel $k$, which gives direct access to $T_k$'s eigendecomposition.

**Theorem 7.** *For an OVK $\mathcal{K} = k_{\mathcal{X}}\,T_k$, an approximate solution to the Huber loss regression problem*

$$\min_{h \in \mathcal{H}_{\mathcal{K}}} \ \frac{1}{n} \sum_{i=1}^{n} \ell_{H,\kappa}(h(x_i) - y_i) + \frac{\Lambda}{2} \|h\|_{\mathcal{H}_{\mathcal{K}}}^2,$$

*is given by Equation* (8), *with* $\hat{\Omega}$ *the solution to the following constrained quadratic problem (with $R$ as in Assumption 6), that can be tackled by PGD in the spirit of Algorithm 1:*

$$\min_{\Omega \in \mathbb{R}^{n \times m}} \quad \mathbf{Tr}\left(\frac{1}{2}\Omega\Omega^\top + \frac{1}{2\Lambda n}K^X\Omega S\Omega^\top - \Omega R^\top\right),$$
$$s.t. \quad \|\Omega\|_{2,\infty} \leqslant \kappa. \tag{9}$$

**Remark 5.** *When $\kappa$ is large, one recovers the unconstrained Ridge regression problem, whose solution enjoys a closed form expression, and for which a resolution method based on an approximation of the inverse of the integral operator $T_k$ was presented in* Kadri et al. (2016).

### 3.3. Stability Analysis

Algorithm stability is a notion introduced by Bousquet and Elisseeff (2002). It links the *stability* of an algorithm, *i.e.* how removing a training observation impacts the algorithm output, to the algorithm generalization capacity, *i.e.* how far the empirical risk of the algorithm output is to its true risk. The rationale behind this approach is that standard analyses of Empirical Risk Minimization rely on a crude approximation consisting in bounding the empirical process $\sup_{h \in \mathcal{H}} |\hat{\mathcal{R}}_n(h) - \mathcal{R}(h)|$. Indeed, considering a supremum over the whole hypothesis set seems very pessimistic, as decision functions with high discrepancy $|\hat{\mathcal{R}}_n(h) - \mathcal{R}(h)|$ would hopefully not be selected by the algorithm. However, the limitation of stability approaches lies in that algorithms performances are never compared to an optimal solution $h^*$. Nevertheless, their capacity to deal with OVK machines without making the trace-class assumption (as opposed to Rademacher-based strategies, see *e.g.* Maurer and Pontil (2016)) make them particularly well suited to our setting. In the footsteps of Audiffren and Kadri (2013), we now derive stability bounds for our algorithms, which are all the more relevant as they make explicit the role of hyperparameters. For any algorithm $A$, $h_{A(\mathcal{S})}$ and $h_{A(\mathcal{S}^{\backslash i})}$ denote the decision functions output by the algorithm, respectively trained on samples $\mathcal{S}$ and $\mathcal{S}^{\backslash i} = \mathcal{S}\backslash\{(x_i, y_i)\}$. Notice that symmetry among observations in Problem (1) cancels the impact of $i$. Formally, algorithm stability states as follows.

**Definition 6.** *(Bousquet and Elisseeff, 2002) Algorithm $A$ has stability $\beta$ if for any sample $\mathcal{S}$, and any $i \leqslant \#\mathcal{S}$, it holds:* $\sup_{(x,y)\in\mathcal{X}\times\mathcal{Y}} |\ell(h_{A(\mathcal{S})}(x), y) - \ell(h_{A(\mathcal{S}^{\backslash i})}(x), y)| \leqslant \beta.$

**Assumption 7.** *There exists $M > 0$ such that for any sample $\mathcal{S}$ and any realization $(x, y) \in \mathcal{X} \times \mathcal{Y}$ of $(X, Y)$ it holds: $\ell(h_{A(\mathcal{S})}(x), y) \leqslant M.$*

**Theorem 8.** *(Bousquet and Elisseeff, 2002) Let $A$ be an algorithm with stability $\beta$ and loss function satisfying Assumption 7. Then, for any $n \geqslant 1$ and $\delta \in ]0, 1[$ it holds with probability at least $1 - \delta$:*

$$\mathcal{R}(h_{A(\mathcal{S})}) \leqslant \hat{\mathcal{R}}_n(h_{A(\mathcal{S})}) + 2\beta + (4n\beta + M)\sqrt{\frac{\ln(1/\delta)}{2n}}.$$

Stability for OVK machines such as in Problem (1) may be derived from the following two assumptions.

**Assumption 8.** *There exists $\gamma > 0$ such that for any input observation $x \in \mathcal{X}$ it holds: $\|\mathcal{K}(x, x)\|_{op} \leqslant \gamma^2.$*

**Assumption 9.** *There exists $C > 0$ such that for any point $(x, y) \in \mathcal{X} \times \mathcal{Y}$, any sample $\mathcal{S}$, and any $i \leqslant \#\mathcal{S}$, it holds: $|\ell(h_{\mathcal{S}}(x), y) - \ell(h_{\mathcal{S}^{\backslash i}}(x), y)| \leqslant C\|h_{\mathcal{S}}(x) - h_{\mathcal{S}^{\backslash i}}(x)\|_{\mathcal{Y}}.$*

**Theorem 9.** *(Audiffren and Kadri, 2013) If Assumptions 8 and 9 hold, then the algorithm returning the solution to Problem (1) has $\beta$ stability with $\beta \leqslant C^2\gamma^2/(\Lambda n).$*

In order to get generalization bounds, we shall now derive constants $M$ and $C$ of Assumptions 7 and 9 respectively. This is usually done under the following assumption.

**Assumption 10.** *There exists $M_{\mathcal{Y}} > 0$ such that for any realization $y \in \mathcal{Y}$ of $Y$ it holds: $\|y\|_{\mathcal{Y}} \leqslant M_{\mathcal{Y}}.$*

**Remark 6.** *It should be noticed that in structured prediction or structured data representation this assumption is directly fulfilled with $M_{\mathcal{Y}} = 1$. Indeed, outputs (and potentially inputs) are actually some $y_i = \phi(z_i)$, with $\phi$ the canonical feature map associated to a scalar kernel, so that it suffices to choose a normalized kernel to satisfy Assumption 10.*

**Theorem 10.** *Under Assumption 10, algorithms previously described satisfy Assumptions 7 and 9 with constants $M$ and $C$ as detailed in Figure 1.*

## 4. Applications and Numerical Experiments

In this section, we discuss some applications unlocked by vv-RKHSs with infinite dimensional outputs. In particular, structured prediction, structured representation learning, and functional regression are formally described, and numerical experiments highlight the benefits of the losses introduced.

### 4.1. Application to Structured Output Prediction

Assume one is interested in learning a predictive decision rule $f$ from a set $\mathcal{X}$ to a complex structured space $\mathcal{Z}$. To bypass the absence of norm on $\mathcal{Z}$, one may design a (scalar) kernel $k$ on $\mathcal{Z}$, whose canonical feature map $\phi : z \mapsto k(\cdot, z)$ transforms any element of $\mathcal{Z}$ into an element of the (scalar) RKHS associated to $k$, denoted $\mathcal{Y}$ ($= \mathcal{H}_k$). Learning a predictive model $f$ from $\mathcal{X}$ to $\mathcal{Z}$ boils down to learning a surrogate vector-valued model $h$ from $\mathcal{X}$ to $\mathcal{Y}$, which is searched for in the vv-RKHS $\mathcal{H}_{\mathcal{K}}$ associated to an OVK $\mathcal{K}$ by solving the following regularized empirical problem.

$$\hat{h} = \underset{h \in \mathcal{H}_{\mathcal{K}}}{\operatorname{argmin}} \ \frac{1}{n}\sum_{i=1}^{n}\ell(h(x_i), \phi(z_i)) + \frac{\Lambda}{2}\|h\|_{\mathcal{H}_{\mathcal{K}}}^2. \tag{10}$$

Once $\hat{h}$ is learned, the predictions in $\mathcal{Z}$ are produced through a pre-image problem $f(x) = \operatorname{argmin}_{z \in \mathcal{Z}} \ell(\phi(z), \hat{h}(x))$. This approach called Input Output Kernel Regression has

| | $M$ | $C$ |
|---|---|---|
| $\epsilon$-SVR | $\sqrt{M_{\mathcal{Y}} - \epsilon}\left(\frac{\sqrt{2}\gamma}{\sqrt{\Lambda}} + \sqrt{M_{\mathcal{Y}} - \epsilon}\right)$ | $1$ |
| $\epsilon$-Ridge | $(M_{\mathcal{Y}} - \epsilon)^2\left(1 + \frac{2\sqrt{2}\gamma}{\sqrt{\Lambda}} + \frac{2\gamma^2}{\Lambda}\right)$ | $2(M_{\mathcal{Y}} - \epsilon)\left(1 + \frac{\gamma\sqrt{2}}{\sqrt{\Lambda}}\right)$ |
| $\kappa$-Huber | $\kappa\sqrt{M_{\mathcal{Y}} - \frac{\kappa}{2}}\left(\frac{\gamma\sqrt{2\kappa}}{\sqrt{\Lambda}} + \sqrt{M_{\mathcal{Y}} - \frac{\kappa}{2}}\right)$ | $\kappa$ |

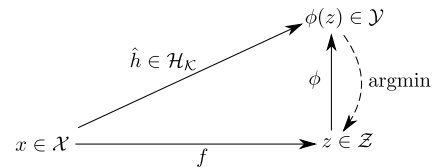*Figure 1.* Algorithms Constants



*Figure 2.* Output Kernel Regression

been studied in several works (Brouard et al., 2011; Kadri et al., 2013). As an instance of the general Output Kernel Regression scheme of Figure 2, it belongs to the family of Surrogate Approaches for structured prediction (see *e.g.* Ciliberto et al. (2016)). While previous works have focused on identity decomposable kernels only, with the squared loss or hinge loss (Brouard et al., 2016b), our general framework allows for many more losses. The use of an $\epsilon$-insensitive loss in Problem (10), in particular, seems adequate as it is a surrogate task, and inducing small mistakes that do not harm the inverse problem, while improving generalization, sounds as a suitable compromise. We thus advocate to solve structured prediction in vv-RKHSs by using losses more sophisticated than the squared norm. In the following, the variants of IOKR are called accordingly to the loss they minimize: $\epsilon$-SV-IOKR, $\epsilon$-Ridge-IOKR, and Huber-IOKR.

**YEAST dataset.** Although our approach's main strength of is to predict infinite dimensional outputs, we start with a simpler standard structured prediction dataset composed of $14$-dimensional outputs (the so-called YEAST dataset Finley and Joachims (2008)) described in the Supplements, on which comparisons and interpretations are easier. We have collected results from Finley and Joachims (2008) and Belanger and McCallum (2016), and benchmarked our three algorithms. Hyperparameters $\Lambda$, $\epsilon$, $\kappa$ have been selected among geometrical grids by cross-validation on the train dataset solely, and performances evaluated on the same test set as the above publications. Results in terms of Hamming error are reported in Figure 6, with significant improvements for the $\epsilon$-Ridge-IOKR and Huber-IOKR. Furthermore, in order to highlight the interactions between our two ways of regularizing, *i.e.* the RKHS norm and the $\epsilon$-insensitivity, we have plotted the $\epsilon$-Ridge-IOKR Mean Square Errors (the Hamming before clamping) and solution sparsity with respect to $\Lambda$ for $\epsilon$ varying from $1e$-5 to 1.5 (Figures 3 and 4): $\Lambda$ and $\epsilon$ seem to act as competitive regularizations. When $\Lambda$ is small, the regularization in $\epsilon$ is efficient, as solution with the best MSE is obtained for $\epsilon$ around $0.6$. Conversely, when $\Lambda$ is big, no sparsity is induced, and having a high $\epsilon$ induces too much regularization. Similar graphs for the $\epsilon$-SVR and $\kappa$-Huber are available in the Supplements, that highlight the superiority of the approaches for a wide range of hyperparameters. A linear output kernel was used, such that solving the inverse problem boils down to clamping.

**Metabolite dataset.** Regarding the infinite dimensional outputs, we have considered the metabolite identification problem (Schymanski et al., 2017), in which one aims at predicting molecules from their mass spectra. For this task, Ridge-IOKR is the state-of-the-art approach, corresponding to our $\epsilon$-Ridge-IOKR with $\epsilon = 0$. Given the high number of constraints, Structured SVMs are not tractable as confirmed by our tests using the Pystruct lib implementation (Müller and Behnke, 2014). This was already noticed in Belanger and McCallum (2016) (14 is the maximum output dimension on which SSVMs were tested), and the implementation we tried indeed yielded very poor results despite prolonged training ($5\%$, $31\%$, $45\%$ top-$k$ errors). We thus investigated the advantages of substituting the standard Ridge Regression for its $\epsilon$-insensitive version or a Huber regression. Outputs (*i.e.* metabolites) are embedded in an infinite dimensional Hilbert space through a Tanimoto-Gaussian kernel with $0.72$ bandwidth. The dataset, presented in the Supplements and described at length in Brouard et al. (2016a), is composed of $6974$ mass spectra, while algorithms are compared through the top-$k$ accuracies, $k = 1, 10, 20$. Two $\Lambda$'s have been picked for their interesting behavior: one that yields the best performance for Ridge-IOKR, and the second that gives the best overall scores (hyperparameters $\epsilon$ and $\kappa$ being chosen to produce the best scores each time). Again, results of Table 1 show improvements due to robust losses that are all the more important as the norm regularization is low, with an improvement on the best overall score.

*Table 1.* Top 1 / 10 / 20 test accuracies (%)

| $\Lambda$ | $1e$-6 | $1e$-4 |
|---|---|---|
| RIDGE-IOKR | 35.7 │ 79.9 │ 86.6 | 38.1 │ 82.0 │ 88.9 |
| $\epsilon$-RIDGE-IOKR | 37.1 │ 81.7 │ 88.3 | 36.3 │ 81.2 │ 87.9 |
| HUBER-IOKR | **38.3** │ **82.2** │ **89.1** | 37.7 │ 81.9 │ 88.8 |

### 4.2. Structured Representation Learning

Extracting vectorial representations from structured inputs is another task that can be tackled in vv-RKHSs (Laforgue et al., 2019). This is a relevant approach in many cases: when complex data are uniquely available under the form of a similarity matrix for instance, for preserving privacy, or when deep neural networks fail to tackle structured objects
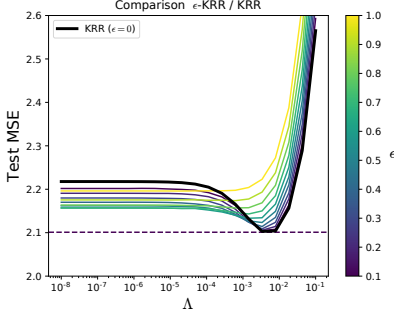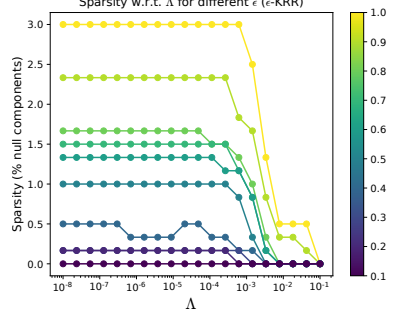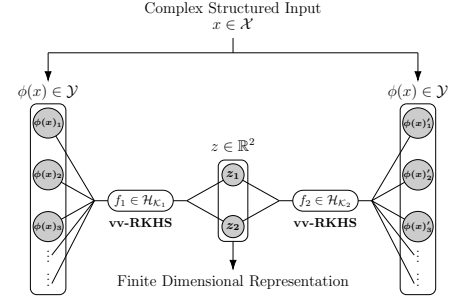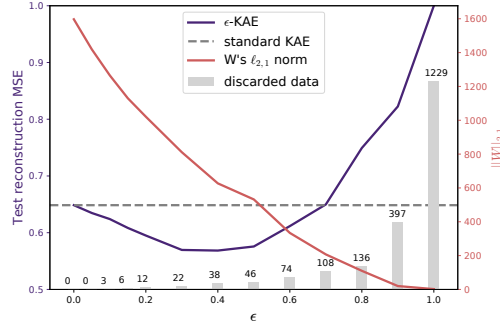
*Figure 3.* Test MSE w.r.t. $\Lambda$



*Figure 4.* Sparsity w.r.t. $\Lambda$



*Figure 5.* 2-Layer Kernel Autoencoder

| | |
|---|---|
| SSVM | 20.2 |
| SPEN | 20.0 |
| $\epsilon$-RIDGE-IOKR | **19.0** |
| HUBER-IOKR | 19.1 |
| $\epsilon$-SV-IOKR | 21.1 |

*Figure 6.* YEAST Hamming errors



*Figure 7.* Reconstruction error w.r.t. $\epsilon$



*Figure 8.* LOO error w.r.t. $\kappa$

as raw data. Embedding data into a Hilbert space makes sense. Then, composing functions in vv-RKHSs results in a Kernel Autoencoder (KAE, Figure 5) that outputs finite codes by minimizing the (regularized) discrepancy:

$$\frac{1}{2n} \sum_{i=1}^{n} \|\phi(x_i) - f_2 \circ f_1(\phi(x_i))\|_{\mathcal{Y}}^2 + \Lambda \operatorname{Reg}(f_1, f_2). \quad (11)$$

Again, this reconstruction loss is not the real goal, but rather a proxy to make the internal representation meaningful. Therefore, all incentives to use $\epsilon$-insensitive losses or the Huber loss still apply. The inferred $\epsilon$-KAE and Huber-KAE, obtained by changing the loss function in Problem (11), are optimized as follows: the first layer coefficients are updated by Gradient Descent, while the second ones are reparametrized into $W_2$ and updated through PGD (instead of KRR closed form for standard KAEs). This has been applied to a drug dataset, introduced in Su et al. (2010) as an extract from the NCI-Cancer database. As shown in Figure 7, the $\epsilon$-insensitivity improves the generalization while inducing sparsity. The $\epsilon$-insensitive framework is thus particularly promising in the context of Autoencoders.

### 4.3. Function-to-Function Regression

Regression with both inputs and outputs of functional nature is a challenging problem at the crossroads of Functional Data Analysis (Ramsay and Silverman, 2007) and Machine Learning (Kadri et al., 2016). While Functional Linear

Modeling is the most common approach to address function-to-function regression, nonparametric approaches based on vv-RKHSs have emerged, that rely on the minimization of a squared loss. However, robustness to abnormal functions is particularly meaningful in a field where data come from sensors and are used to monitor physical assets. To the best of our knowledge, robust regression has only been tackled in the context of Functional Linear Models (Kalogridis and Van Aelst, 2019). We propose here to highlight the relevance of OVK machines learned with a Huber loss by solving Problem (9) for various levels $\kappa$.

**Lip acceleration from EMG dataset.** We consider the problem of predicting lip acceleration among time from electromyography (EMG) signals (Ramsay and Silverman, 2007). The dataset consists of 32 records of the lower lip trajectory over 641 timestamps, and the associated EMG records, augmented with 4 outliers to assess the robustness of our approach. Usefulness of minimizing the Huber loss is illustrated in Figure 8 by computing the Leave-One-Out (LOO) error associated to each model for various values of $m$. For each $m$, as $\kappa$ grows larger than a threshold, the constraint on $\|\Omega\|_{2,\infty}$ becomes void and we recover the Ridge Regression solution. The kernel chosen is given by $k_{\mathcal{X}}(x_1, x_2) = \int_0^1 \exp\left(|x_1(\theta) - x_2(\theta)|\right) \mathrm{d}\theta$, with $(\psi_j)_{j=1}^m$ being the harmonic basis in sine and cosine of $L^2[0, 1]$, and $(\lambda_j)_{j=1}^m = (1/(j+1)^2)_{j=1}^m$.

## 4.4. Related Work

Another application of the presented results, both theoretical and computational, is the generalization of the *loss trick*, see *e.g.* Ciliberto et al. (2016). In the context of Output Kernel Regression, the latter stipulates that for suitable losses, the decoding expresses in terms of loss evaluations. The work by Luise et al. (2019) has extended this trick to penalization schemes different from the natural vv-RKHS norm. Our findings, and the double expansion in particular, suggest that the loss trick can still be used with other surrogate loss functions than the squared norm, opening the door to a wide range of applications.

## 5. Conclusion

This work presents a versatile framework based on duality to learn OVK machines with infinite dimensional outputs. The case of convolved losses (*e.g.* $\epsilon$-insensitive, Huber) is thoroughly tackled, from algorithmic procedures to stability analysis. This offers novel ways to enforce sparsity and robustness when learning within vv-RKHSs, opening an avenue for new applications on structured and functional data (*e.g.* anomaly detection, robust prediction). Future research directions could feature a calibration study of these novel surrogate approaches, or the introduction of kernel approximations such as random Fourier features, that would benefit our framework twice: both in input and in output.

REFERENCES

Álvarez, M. A., Rosasco, L., and Lawrence, N. D. (2012). Kernels for vector-valued functions: a review. *Foundations and Trends in Machine Learning*, 4(3):195–266.

Audiffren, J. and Kadri, H. (2013). Stability of multi-task kernel regression algorithms. In *Asian Conference on Machine Learning*, pages 1–16.

Baldassarre, L., Rosasco, L., Barla, A., and Verri, A. (2012). Multi-output learning via spectral filtering. *Machine Learning*, 87(3):259–301.

Bauschke, H. H., Combettes, P. L., et al. (2011). *Convex analysis and monotone operator theory in Hilbert spaces*, volume 408. Springer.

Belanger, D. and McCallum, A. (2016). Structured prediction energy networks. In *International Conference on Machine Learning*, pages 983–992.

Bousquet, O. and Elisseeff, A. (2002). Stability and generalization. *Journal of machine learning research*, 2(Mar):499–526.

Brault, R., Lambert, A., Szabo, Z., Sangnier, M., and d'Alché-Buc, F. (2019). Infinite task learning in rkhss. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pages 1294–1302.

Brouard, C., d'Alché-Buc, F., and Szafranski, M. (2011). Semi-supervised penalized output kernel regression for link prediction. In *International Conference on Machine Learning (ICML)*, pages 593–600.

Brouard, C., Shen, H., Dührkop, K., d'Alché-Buc, F., Böcker, S., and Rousu, J. (2016a). Fast metabolite identification with input output kernel regression. *Bioinformatics*, 32(12):28–36.

Brouard, C., Szafranski, M., and D'Alché-Buc, F. (2016b). Input output kernel regression: supervised and semi-supervised structured output prediction with operator-valued kernels. *The Journal of Machine Learning Research*, 17(1):6105–6152.

Caponnetto, A., Micchelli, C. A., Pontil, M., and Ying, Y. (2008). Universal multitask kernels. *Journal of Machine Learning Research*, 9:1615–1646.

Carmeli, C., De Vito, E., and Toigo, A. (2006). Vector valued reproducing kernel hilbert spaces of integrable functions and mercer theorem. *Analysis and Applications*, 4(04):377–408.

Carmeli, C., De Vito, E., Toigo, A., and Umanitá, V. (2010). Vector valued reproducing kernel hilbert spaces and universality. *Analysis and Applications*, 8(01):19–61.

Ciliberto, C., Rosasco, L., and Rudi, A. (2016). A consistent regularization approach for structured prediction. In *Advances in Neural Information Processing Systems (NIPS) 29*, pages 4412–4420.

Dinuzzo, F., Ong, C., Gehler, P., and Pillonetto, G. (2011). Learning output kernels with block coordinate descent. In *International Conference on Machine Learning (ICML)*, pages 49–56.

Drucker, H., Burges, C. J., Kaufman, L., Smola, A. J., and Vapnik, V. (1997). Support vector regression machines. In *Advances in neural information processing systems*, pages 155–161.

Finley, T. and Joachims, T. (2008). Training structural svms when exact inference is intractable. In *Proceedings of the 25th international conference on Machine learning*, pages 304–311.

Fung, G. and Mangasarian, O. L. (2000). Data selection for support vector machine classifiers. In Ramakrishnan, R., Stolfo, S. J., Bayardo, R. J., and Parsa, I., editors, *Proceedings of the sixth ACM SIGKDD international conference on Knowledge discovery and data mining, Boston, MA, USA, August 20-23, 2000*, pages 64–70. ACM.

Hastie, T., Tibshirani, R., and Wainwright, M. (2015). *Statistical learning with sparsity: the lasso and generalizations*. Chapman and Hall/CRC.

Hofmann, T., Schoelkopf, B., and Smola, A. J. (2008). Kernel methods in machine learning. *Ann. Statist.*, 36(3):1171–1220.

Huber, P. J. (1964). Robust estimation of a location parameter. *The Annals of Mathematical Statistics*, pages 73–101.

Joachims, T., Hofmann, T., Yue, Y., and Yu, C.-N. (2009). Predicting structured objects with support vector machines. *Commun. ACM*, 52(11):97–104.

Kadri, H., Duflos, E., Preux, P., Canu, S., and Davy, M. (2010). Nonlinear functional regression: a functional rkhs approach. In Teh, Y. W. and Titterington, M., editors, *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics*, volume 9 of *Proceedings of Machine Learning Research*, pages 374–380. PMLR.

Kadri, H., Duflos, E., Preux, P., Canu, S., Rakotomamonjy, A., and Audiffren, J. (2016). Operator-valued kernels for learning from functional response data. *Journal of Machine Learning Research*, 17(20):1–54.

Kadri, H., Ghavamzadeh, M., and Preux, P. (2013). A generalized kernel approach to structured output learning. In *International Conference on Machine Learning (ICML)*, pages 471–479.

Kalogridis, I. and Van Aelst, S. (2019). Robust functional regression based on principal components. *Journal of Multivariate Analysis*, 173:393 – 415.

Laforgue, P., Clémençon, S., and d'Alché-Buc, F. (2019). Autoencoding any data through kernel autoencoders. In *Artificial Intelligence and Statistics*, pages 1061–1069.

Luise, G., Stamos, D., Pontil, M., and Ciliberto, C. (2019). Leveraging low-rank relations between surrogate tasks in structured prediction. *arXiv preprint arXiv:1903.00667*.

Maurer, A. and Pontil, M. (2016). Bounds for vector-valued function estimation. *arXiv preprint arXiv:1606.01487*.

Micchelli, C. A. and Pontil, M. (2005). On learning vector-valued functions. *Neural computation*, 17(1):177–204.

Moreau, J. J. (1962). Fonctions convexes duales et points proximaux dans un espace hilbertien. *Comptes rendus hebdomadaires des séances de l'Académie des sciences*, 255:2897–2899.

Müller, A. C. and Behnke, S. (2014). pystruct - learning structured prediction in python. *Journal of Machine Learning Research*, 15:2055–2060.

Obozinski, G., Taskar, B., and Jordan, M. I. (2010). Joint covariate selection and joint subspace selection for multiple classification problems. *Statistics and Computing*, 20(2):231–252.

Osborn, J. E. (1975). Spectral approximation for compact operators. *Mathematics of computation*, 29(131):712–725.

Ramsay, J. O. and Silverman, B. W. (2007). *Applied functional data analysis: methods and case studies*. Springer.

Sangnier, M., Fercoq, O., and d'Alché-Buc, F. (2017). Data sparse nonparametric regression with $\epsilon$-insensitive losses. In *Asian Conference on Machine Learning*, pages 192–207.

Schymanski, E., Ruttkies, C., and Krauss, M. e. a. (2017). Critical assessment of small molecule identification 2016: automated methods. *J Cheminform*, 9:22.

Senkene, E. and Tempel'man, A. (1973). Hilbert spaces of operator-valued functions. *Lithuanian Mathematical Journal*, 13(4):665–670.

Su, H., Heinonen, M., and Rousu, J. (2010). Structured output prediction of anti-cancer drug activity. In *IAPR International Conference on Pattern Recognition in Bioinformatics*, pages 38–49. Springer.

Tsochantaridis, I., Joachims, T., Hofmann, T., and Altun, Y. (2005). Large margin methods for structured and interdependent output variables. *Journal of machine learning research*, 6(Sep):1453–1484.

Zhu, J., Hoi, S. C. H., and Lyu, M. R. (2008). Robust regularized kernel regression. *IEEE Transactions on Systems, Man, and Cybernetics, Part B (Cybernetics)*, 38(6):1639–1644.

The Supplementary Material is organized as follows. Appendix A collects the technical proofs of the core article's results. Appendix B provides illustrations of the main loss functions considered ($\epsilon$-insensitive Ridge and SVR, $\kappa$-Huber) in 1 and 2 dimensions. Appendix C gathers additional details about the experimental protocols and the code furnished.

## A. Technical Proofs

### A.1. Proof of Theorem 3

First, notice that the primal problem

$$\min_{h \in \mathcal{H}_{\mathcal{K}}} \quad \frac{1}{n} \sum_{i=1}^{n} \ell(h(x_i), y_i) + \frac{\Lambda}{2} \|h\|_{\mathcal{H}_{\mathcal{K}}}^2$$

can be rewritten

$$\min_{h \in \mathcal{H}_{\mathcal{K}}} \quad \sum_{i=1}^{n} \ell_i(u_i) + \frac{\Lambda n}{2} \|h\|_{\mathcal{H}_{\mathcal{K}}}^2,$$

$$\text{s.t.} \quad u_i = h(x_i) \quad \forall i \leqslant n.$$

Therefore, with the notation $\boldsymbol{u} = (u_i)_{i \leqslant n}$ and $\boldsymbol{\alpha} = (\alpha_i)_{i \leqslant n}$, the Lagrangian writes

$$\mathscr{L}(h, \boldsymbol{u}, \boldsymbol{\alpha}) = \sum_{i=1}^{n} \ell_i(u_i) + \frac{\Lambda n}{2} \|h\|_{\mathcal{H}_{\mathcal{K}}}^2 + \sum_{i=1}^{n} \langle \alpha_i, u_i - h(x_i) \rangle_{\mathcal{Y}},$$

$$= \sum_{i=1}^{n} \ell_i(u_i) + \frac{\Lambda n}{2} \|h\|_{\mathcal{H}_{\mathcal{K}}}^2 + \sum_{i=1}^{n} \langle \alpha_i, u_i \rangle_{\mathcal{Y}} - \sum_{i=1}^{n} \langle \mathcal{K}(\cdot, x_i)\alpha_i, h \rangle_{\mathcal{H}_{\mathcal{K}}}.$$

Differentiating with respect to $h$ and using the definition of the Fenchel-Legendre transform, one gets

$$g(\boldsymbol{\alpha}) = \inf_{h \in \mathcal{H}_{\mathcal{K}}, \boldsymbol{u} \in \mathcal{Y}^n} \mathscr{L}(h, \boldsymbol{u}, \boldsymbol{\alpha}),$$

$$= \sum_{i=1}^{n} \inf_{u_i \in \mathcal{Y}} \left\{ \ell_i(u_i) + \langle \alpha_i, u_i \rangle_{\mathcal{Y}} \right\} + \inf_{h \in \mathcal{H}_{\mathcal{K}}} \left\{ \frac{\Lambda n}{2} \|h\|_{\mathcal{H}_{\mathcal{K}}}^2 - \sum_{i=1}^{n} \langle \mathcal{K}(\cdot, x_i)\alpha_i, h \rangle_{\mathcal{H}_{\mathcal{K}}} \right\},$$

$$= \sum_{i=1}^{n} -\ell_i^{\star}(-\alpha_i) - \frac{1}{2\Lambda n} \sum_{i,j=1}^{n} \langle \alpha_i, \mathcal{K}(x_i, x_j)\alpha_j \rangle_{\mathcal{Y}},$$

together with the equality $\hat{h} = \frac{1}{\Lambda n} \sum_{i=1}^{n} \mathcal{K}(\cdot, x_i)\alpha_i$. The conclusion follows immediately. □

### A.2. Proof of Theorem 4

As a reminder, our goal is to compute the solutions to the following problem:

$$\hat{h} \in \underset{h \in \mathcal{H}_{\mathcal{K}}}{\operatorname{argmin}} \quad \frac{1}{n} \sum_{i=1}^{n} \ell(h(x_i), y_i) + \frac{\Lambda}{2} \|h\|_{\mathcal{H}_{\mathcal{K}}}^2.$$

Using Theorem 3, one gets that $\hat{h} = \frac{1}{\Lambda n} \sum_{i=1}^{n} \mathcal{K}(\cdot, x_i)\hat{\alpha}_i$, with the $(\hat{\alpha}_i)_{i \leqslant n}$ satisfying:

$$(\hat{\alpha}_i)_{i=1}^{n} \in \underset{(\alpha_i)_{i=1}^{n} \in \mathcal{Y}^n}{\operatorname{argmin}} \quad \sum_{i=1}^{n} \ell_i^{\star}(-\alpha_i) + \frac{1}{2\Lambda n} \sum_{i,j=1}^{n} \langle \alpha_i, \mathcal{K}(x_i, x_j)\alpha_j \rangle_{\mathcal{Y}}.$$

However, this optimization problem cannot be solved in a straightforward manner, as $\mathcal{Y}$ is in general infinite dimensional. Nevertheless, it is possible to bypass this difficulty by noticing that the optimal $(\hat{\alpha}_i)_{i \leqslant n}$ actually lie in $\mathbf{Y}^n$. To show this, we decompose each coefficient as $\hat{\alpha}_i = \alpha_i^{\mathbf{Y}} + \alpha_i^{\perp}$, with $(\alpha_i^{\mathbf{Y}})_{i \leqslant n}, (\alpha_i^{\perp})_{i \leqslant n} \in \mathbf{Y}^n \times \mathbf{Y}^{\perp n}$. Then, noticing that non-null $(\alpha_i^{\perp})_{i \leqslant n}$ necessarily increase the objective, we can conclude that the optimal $(\hat{\alpha}_i)_{i \leqslant n}$ have no components among $\mathbf{Y}^{\perp}$, or equivalently pertain to $\mathbf{Y}$. Indeed, by virtue of Assumptions 1 and 3, it holds:

$$\sum_{i=1}^{n} \ell_i^{\star}(-\alpha_i^{\mathbf{Y}}) + \frac{1}{2\Lambda n} \sum_{i,j=1}^{n} \langle \alpha_i^{\mathbf{Y}}, \mathcal{K}(x_i, x_j)\alpha_j^{\mathbf{Y}} \rangle_{\mathcal{Y}} \leqslant \sum_{i=1}^{n} \ell_i^{\star}(-\alpha_i^{\mathbf{Y}} - \alpha_i^{\perp}) + \frac{1}{2\Lambda n} \sum_{i,j=1}^{n} \langle \alpha_i^{\mathbf{Y}} + \alpha_i^{\perp}, \mathcal{K}(x_i, x_j)(\alpha_j^{\mathbf{Y}} + \alpha_j^{\perp}) \rangle_{\mathcal{Y}}.$$

If the inequality about $\ell_i^{\star}$ follows directly Assumption 1, that about $\mathcal{K}(x_i, x_j)$ can be obtained by Assumption 3 as follows:

$$\sum_{i,j=1}^{n} \langle \alpha_i^{\mathbf{Y}} + \alpha_i^{\perp}, \mathcal{K}(x_i, x_j)(\alpha_j^{\mathbf{Y}} + \alpha_j^{\perp}) \rangle_{\mathcal{Y}}$$

$$= \sum_{i,j=1}^{n} \langle \alpha_i^{\mathbf{Y}}, \mathcal{K}(x_i, x_j)\alpha_j^{\mathbf{Y}} ) \rangle_{\mathcal{Y}} + 2 \sum_{i,j=1}^{n} \langle \alpha_i^{\perp}, \mathcal{K}(x_i, x_j)\alpha_j^{\mathbf{Y}} \rangle_{\mathcal{Y}} + \sum_{i,j=1}^{n} \langle \alpha_i^{\perp}, \mathcal{K}(x_i, x_j)\alpha_j^{\perp} \rangle_{\mathcal{Y}},$$

$$= \sum_{i,j=1}^{n} \langle \alpha_i^{\mathbf{Y}}, \mathcal{K}(x_i, x_j)\alpha_j^{\mathbf{Y}} ) \rangle_{\mathcal{Y}} + \sum_{i,j=1}^{n} \langle \alpha_i^{\perp}, \mathcal{K}(x_i, x_j)\alpha_j^{\perp} \rangle_{\mathcal{Y}},$$

$$\geqslant \sum_{i,j=1}^{n} \langle \alpha_i^{\mathbf{Y}}, \mathcal{K}(x_i, x_j)\alpha_j^{\mathbf{Y}} ) \rangle_{\mathcal{Y}},$$

where we have used successively Assumption 3 and the positiveness of $\mathcal{K}$. So there exists $\Omega = [\omega_{ij}]_{1 \leqslant i,j \leqslant n} \in \mathbb{R}^{n \times n}$ such that for all $i \leqslant n$, $\hat{\alpha}_i = \sum_j \omega_{ij} y_j$. This proof technique is very similar in spirit to that of the Representer Theorem, and yields an analogous result, the reduction of the search space to a smaller vector space, as discussed at length in the main text. The dual optimization problem thus rewrites:

$$\min_{\Omega \in \mathbb{R}^{n \times n}} \sum_{i=1}^{n} \ell_i^{\star}\left( -\sum_{j=1}^{n} \omega_{ij} \, y_j \right) + \frac{1}{2\Lambda n} \sum_{i,j=1}^{n} \left\langle \sum_{k=1}^{n} \omega_{ik} \, y_k, \mathcal{K}(x_i, x_j) \sum_{l=1}^{n} \omega_{jl} \, y_l \right\rangle_{\mathcal{Y}}$$

$$\min_{\Omega \in \mathbb{R}^{n \times n}} \sum_{i=1}^{n} L_i\left( (\omega_{ij})_{j \leqslant n}, K^Y \right) + \frac{1}{2\Lambda n} \sum_{i,j,k,l=1}^{n} \omega_{ik} \, \omega_{jl} \left\langle y_k, \sum_{t=1}^{T} k_t(x_i, x_j)A_t y_l \right\rangle_{\mathcal{Y}},$$

$$\min_{\Omega \in \mathbb{R}^{n \times n}} \sum_{i=1}^{n} L_i\left( \Omega_{i:}, K^Y \right) + \frac{1}{2\Lambda n} \mathbf{Tr}\left( \tilde{M}^{\top}(\Omega \otimes \Omega) \right), \tag{12}$$

with $M$ the $n \times n \times n \times n$ tensor such that $M_{ijkl} = \langle y_k, \mathcal{K}(x_i, x_j)y_l \rangle_{\mathcal{Y}}$, and $\tilde{M}$ its rewriting as a $n^2 \times n^2$ block matrix such that its $(i, j)$ block is the $n \times n$ matrix with elements $\tilde{M}_{st}^{(i,j)} = \langle y_j, \mathcal{K}(x_i, x_s)y_t \rangle_{\mathcal{Y}}$.

The second term is quadratic in $\Omega$, and consequently convex. As for the $L_i$'s, they are basically rewritings of the Fenchel-Legendre transforms $\ell_i^{\star}$'s that ensure the computability of the problem (they only depend on $K^Y$, which is known). Regarding their convexity, we know by definition that the $\ell_i^{\star}$'s are convex. Composing by a linear function preserving the convexity, we know that each $L_i$ is convex with respect to $\Omega_{i:}$, and therefore with respect to $\Omega$.

Thus, we have first converted the infinite dimensional primal problem in $\mathcal{H}_{\mathcal{K}}$ into an infinite dimensional dual problem in $\mathcal{Y}^n$, which in turn is reduced to a convex optimization procedure over $\mathbb{R}^{n \times n}$, that only involves computable quantities.

If $\mathcal{K}$ satisfies Assumption 4, the tensor $M$ simplifies to

$$M_{ijkl} = \langle y_k, \mathcal{K}(x_i, x_j)y_l \rangle_{\mathcal{Y}} = \sum_{t=1}^{T} k_t(x_i, x_j) \langle y_k, A_t y_l \rangle_{\mathcal{Y}} = \sum_{t=1}^{T} [K_t^X]_{ij}[K_t^Y]_{kl},$$

and the problem rewrites

$$\min_{\Omega \in \mathbb{R}^{n \times n}} \sum_{i=1}^{n} L_i\left( \Omega_{i:}, K^Y \right) + \frac{1}{2\Lambda n} \sum_{t=1}^{T} \mathbf{Tr}\left( K_t^X \Omega K_t^Y \Omega^{\top} \right).$$

$\square$

**Remark 7.** *The second term of Problem* (12) *can be easily optimized. Indeed, let $\tilde{M}$ be a block matrix such that $\tilde{M}_{st}^{(i,j)} = \tilde{M}_{ij}^{(s,t)}$ for all $i, j, s, t \leqslant n$. Notice that $\tilde{M}$ as defined earlier satisfies this condition as a direct consequence of the OVK symmetry property. Then it holds:*

$$\frac{\partial \mathbf{Tr}\left(\tilde{M}^{\top}(\Omega \otimes \Omega)\right)}{\partial \omega_{st}} = 2\mathbf{Tr}\left(\tilde{M}^{(s,t)\top}\Omega\right).$$

Indeed, notice that $\mathbf{Tr}\left(\tilde{M}^{\top}(\Omega \otimes \Omega)\right) = \sum_{i,j=1}^{n} \omega_{ij}\mathbf{Tr}\left(\tilde{M}^{(i,j)\top}\Omega\right)$ and use the symmetry assumption. In the particular case of a decomposable kernel, it holds that $\tilde{M}^{(i,j)} = K_{i:}^{X}K_{j:}^{Y\top}$ so that

$$\frac{\partial \mathbf{Tr}\left(\tilde{M}^{\top}(\Omega \otimes \Omega)\right)}{\partial \omega_{st}} = 2\mathbf{Tr}\left(\tilde{M}^{(s,t)\top}\Omega\right) = 2\sum_{i,j=1}^{n}\left[K_{s:}^{X}K_{t:}^{Y\top}\right]_{ij}\omega_{ij} = 2\sum_{ij=1}^{n}K_{si}^{X}K_{tj}^{Y}\omega_{ij} = 2\left[K^{X}\Omega K^{Y}\right]_{st},$$

and one recovers the gradients established in Equation (15).

## A.3. Proof of Proposition 1

The proof technique is the same for all losses: first explicit the FL transforms $\ell_i^{\star}$, then use simple arguments to verify Assumptions 1 and 2. For instance, any increasing function of $\|\alpha\|$ automatically satisfy the assumptions.

- Assume that $\ell$ is such that there is $f : \mathbb{R} \to \mathbb{R}$ convex, $\forall i \leqslant n, \exists z_i \in Y, \ell_i(y) = f(\langle y, z_i \rangle)$. Then $\ell_i^{\star} : \mathcal{Y} \to \mathbb{R}$ writes $\ell_i^{\star}(\alpha) = \sup_{y \in \mathcal{Y}} \langle \alpha, y \rangle - f(\langle y, z_i \rangle)$. If $\alpha$ is not collinear to $z_i$, this quantity is obviously $+\infty$. Otherwise, assume that $\alpha = \lambda z_i$. The FL transform rewrites: $\ell_i^{\star}(\alpha) = \sup_t \lambda t - f(t) = f^{\star}(\lambda) = f^{\star}(\pm\|\alpha\|/\|z_i\|)$. Finally, $\ell_i^{\star}(\alpha) = \chi_{\mathrm{span}(z_i)}(\alpha) + f^{\star}\left(\pm\frac{\|\alpha\|}{\|z_i\|}\right)$. If $\alpha \notin Y$, then *a fortiori* $\alpha \notin \mathrm{span}(z_i)$, so $\ell_i^{\star}(\alpha^Y + \alpha^{\perp}) = +\infty \geqslant \ell_i^{\star}(\alpha^Y)$ for all $(\alpha^Y, \alpha^{\perp}) \in Y \times Y^{\perp}$. For all $i \leqslant n$, $\ell_i^{\star}$ satisfy Assumption 1. As for Assumption 2, if $\alpha = \sum_{i=1}^{n} c_i y_i$, then $\chi_{\mathrm{span}(z_i)}(\alpha)$ only depends on the $(c_i)_{i \leqslant n}$ Indeed, assume that $z_i \in Y$ writes $\sum_j b_j y_j$. Then $\chi_{\mathrm{span}(z_i)}(\alpha)$ is equal to 0 if there exists $\lambda \in \mathbb{R}$ such that $c_j = \lambda b_j$ for all $j \leqslant n$, and to $+\infty$ otherwise. The second term of $\ell_i^{\star}$ depending only on $\|\alpha\|$, it directly satisfies Assumption 2. This concludes the proof.

- Assume that $\ell$ is such that there is $f : \mathbb{R}_+ \to \mathbb{R}$ convex increasing, with $\frac{f'(t)}{t}$ continuous over $\mathbb{R}_+$, $\ell(y) = f(\|y\|)$. Although this loss may seem useless at the first sight since $\ell$ does not depend on $y_i$, it should not be forgotten that the composition with $y \mapsto y - y_i$ does not affect the validation of Assumptions 1 and 2 (see below). One has: $\ell^{\star}(\alpha) = \sup_{y \in \mathcal{Y}} \langle \alpha, y \rangle - f(\|y\|)$. Differentiating w.r.t. $y$, one gets: $\alpha = \frac{f'(\|y\|)}{\|y\|}y$, which is always well define as $t \mapsto \frac{f'(t)}{t}$ is continuous over $\mathbb{R}_+$. Reverting the equality, it holds: $y = \frac{f'^{-1}(\|\alpha\|)}{\|\alpha\|}\alpha$, and $\ell^{\star}(\alpha) = \|\alpha\|f'^{-1}(\|\alpha\|) - f \circ f'^{-1}(\|\alpha\|)$. This expression depending only on $\|\alpha\|$, Assumption 2 is automatically satisfied. Let us now investigate the monotonicity of $\ell^{\star}$ w.r.t. $\|\alpha\|$. Let $g : \mathbb{R}_+ \to \mathbb{R}$ such that $g(t) = tf'^{-1}(t) - f \circ f'^{-1}(t)$. Then $g'(t) = f'^{-1}(t) \geqslant 0$. Indeed, as $f' : \mathbb{R}_+ \to \mathbb{R}_+$ is always positive due to the monotonicity of $f$, so is $f'^{-1}$. This final remark guarantees that $\ell^{\star}$ is increasing with $\|\alpha\|$. It is then direct that $\ell^{\star}$ fulfills Assumption 1.

- Assume that $\ell(y) = \lambda\|y\|$. It holds $\ell^{\star}(\alpha) = \chi_{\mathcal{B}_{\lambda}}(\alpha)$. So $\ell^{\star}$ is increasing w.r.t. $\|\alpha\|$: it fulfills Assumptions 1 and 2.

- Assume that $\ell(y) = \chi_{\mathcal{B}_{\lambda}}(y)$. It holds $\ell^{\star}(\alpha) = \lambda\|\alpha\|$. The monotonicity argument also applies.

- Assume that $\ell(y) = \lambda\|y\|\log(\|y\|)$. It can be shown that $\ell^{\star}(\alpha) = \lambda e^{\frac{\|\alpha\|}{\lambda} - 1}$. The same argument as above applies.

- Assume that $\ell(y) = \lambda(\exp(\|y\|) - 1)$. It can be shown that $\ell^{\star}(\alpha) = \mathbb{I}\{\|\alpha\| \geqslant \lambda\} \cdot \left(\|\alpha\|\log\left(\frac{\|\alpha\|}{\lambda e}\right) + \lambda\right)$. Again, the FL transform is an increasing function of $\|\alpha\|$: it satisfies Assumptions 1 and 2.

- Assume that $\ell_i(y) = f(y - y_i)$, with f such that $f^{\star}$ fulfills Assumptions 1 and 2. Then $\ell_i^{\star}(\alpha) = \sup_{y \in \mathcal{Y}} \langle \alpha, y \rangle - f(y - y_i) = f^{\star}(\alpha) + \langle \alpha, y_i \rangle$. If $f^{\star}$ satisfies Assumptions 1 and 2, then so does $\ell_i^{\star}$. This remark is very important, as it gives more sense to loss function based on $\|y\|$ only, since they can be applied to $y - y_i$ now.

- Assume that there exists $f, g$ satisfying Assumptions 1 and 2 such that $\ell_i(y) = (f \square g)(y)$, where $\square$ denotes the infimal convolution, *i.e.* $(f \square g)(y) = \inf_x f(x) + g(y - x)$. Standard arguments about FL transforms state that $(f \square g)^{\star} = f^{\star} + g^{\star}$, so that if both $f$ and $g$ satisfy Assumptions 1 and 2, so does $f \square g$. This last example allows to deal with $\epsilon$-insensitive losses for instance (convolution of a loss and $\chi_{\mathcal{B}_{\epsilon}}$), the Huber loss (convolution of $\|.\|$ and $\|.\|^2$), or more generally all Moreau envelopes (convolution of a loss and $\frac{1}{2}\|.\|^2$).

$\square$

## A.4. Proof of Theorem 5

The proof of Theorem 5 is straightforward: since the dual space $\tilde{\mathcal{Y}}_m$ is of finite dimension $m$, the dual variable can be written as a linear combination of the $\{\psi_j\}_{j=1}^m$ to get Problem (7).

## A.5. Proof of Theorem 6

### A.5.1. $\epsilon$-RIDGE – FROM PROBLEM $(P1)$ TO $(D1)$

Applying Theorem 3 together with the Fenchel-Legendre transforms detailed in the proof of Proposition 1, a dual to the $\epsilon$-Ridge regression primal problem is:

$$\min_{(\alpha_i)_{i=1}^n \in \mathcal{Y}^n} \quad \frac{1}{2}\sum_{i=1}^n \|\alpha_i\|_{\mathcal{Y}}^2 - \sum_{i=1}^n \langle \alpha_i, y_i \rangle_{\mathcal{Y}} + \epsilon \sum_{i=1}^n \|\alpha_i\|_{\mathcal{Y}} + \frac{1}{2\Lambda n}\sum_{i,j=1}^n \langle \alpha_i, \mathcal{K}(x_i, x_j)\alpha_j \rangle_{\mathcal{Y}},$$

$$\min_{(\alpha_i)_{i=1}^n \in \mathcal{Y}^n} \quad \frac{1}{2}\sum_{i,j=1}^n \left\langle \alpha_i, \left(\delta_{ij}\mathbf{I}_{\mathcal{Y}} + \frac{1}{\Lambda n}\mathcal{K}(x_i, x_j)\right)\alpha_j \right\rangle_{\mathcal{Y}} - \sum_{i=1}^n \langle \alpha_i, y_i \rangle_{\mathcal{Y}} + \epsilon \sum_{i=1}^n \|\alpha_i\|_{\mathcal{Y}}.$$

By virtue of Theorem 4, we known that the optimal $(\alpha_i)_{i=1}^n \in \mathcal{Y}^n$ are in $\mathbf{Y}^n$. After the reparametrization $\alpha_i = \sum_j \omega_{ij}\, y_j$, the problem rewrites:

$$\min_{\Omega \in \mathbb{R}^{n \times n}} \quad \frac{1}{2}\mathbf{Tr}\left(\tilde{K}\Omega K^Y \Omega^\top\right) - \mathbf{Tr}\left(K^Y \Omega\right) + \epsilon \sum_{i=1}^n \sqrt{[\Omega K^Y \Omega^\top]_{ii}}, \tag{13}$$

with $\Omega$, $\tilde{K}$, $K^Y$ the $n \times n$ matrices such that $[\Omega]_{ij} = \omega_{ij}$, $\tilde{K} = \frac{1}{\Lambda n}K^X + \mathbf{I}_n$, and $[K^Y]_{ij} = \langle y_i, y_j \rangle_{\mathcal{Y}}$.

Now, let $K^Y = U\Sigma U^\top = \left(U\Sigma^{1/2}\right)\left(U\Sigma^{1/2}\right)^\top = VV^\top$ be the SVD of $K^Y$, and let $W = \Omega V$. Notice that $K^Y$ is positive semi-definite, and can be made positive definite if necessary, so that $V$ is full rank, and optimizing with respect to $W$ is strictly equivalent to minimizing with respect to $\Omega$. With this change of variable, Problem (13) rewrites:

$$\min_{W \in \mathbb{R}^{n \times n}} \quad \frac{1}{2}\mathbf{Tr}\left(\tilde{K}WW^\top\right) - \mathbf{Tr}\left(V^\top W\right) + \epsilon\|W\|_{2,1}, \tag{14}$$

with $\|W\|_{2,1} = \sum_i \|W_{i:}\|_2$ the row-wise $\ell_{2,1}$ mixed norm of matrix $W$. With $\tilde{K} = A^\top A$ the SVD of $\tilde{K}$, and $B$ such that $A^\top B = V$, one can add the constant term $\frac{1}{2}\mathbf{Tr}(A^{\top^{-1}}VV^\top A^{-1}) = \frac{1}{2}\mathbf{Tr}(BB^\top)$ to the objective without changing Problem (14). One finally gets the Multi-Task Lasso problem:

$$\min_{W \in \mathbb{R}^{n \times n}} \quad \frac{1}{2}\|AW - B\|_{\mathrm{Fro}}^2 + \epsilon\|W\|_{2,1}.$$

$\square$

We also emphasize that we recover the solution to the standard Ridge regression when $\epsilon = 0$. Indeed, coming back to Problem (13) and differentiating with respect to $\Omega$, one gets:

$$\tilde{K}\hat{\Omega}K^Y - K^Y = 0 \iff \hat{\Omega} = \tilde{K}^{-1},$$

which is exactly the standard kernel Ridge regression solution, see *e.g.* Brouard et al. (2016b).

Furthermore, notice that when $\mathcal{K}$ is not identity decomposable, but only satisfies Assumption 4, then Problem (14) cannot be factorized that easily. Nonetheless, it admits a simple resolution, as detailed in the following lines. After the $\Omega$ reparametrization, the problem writes

$$\min_{\Omega \in \mathbb{R}^{n \times n}} \quad \frac{1}{2}\mathbf{Tr}(\Omega K^Y \Omega^\top) - \mathbf{Tr}\left(K^Y \Omega\right) + \epsilon \sum_{i=1}^n \sqrt{[\Omega K^Y \Omega^\top]_{i,i}} + \frac{1}{2\Lambda n}\sum_{t=1}^T \mathbf{Tr}(K_t^X \Omega K_t^Y \Omega^\top),$$

$$\min_{W \in \mathbb{R}^{n \times n}} \quad \frac{1}{2}\mathbf{Tr}(WW^\top) + \frac{1}{2\Lambda n}\sum_{t=1}^T \mathbf{Tr}(K_t^X W \tilde{K}_t^Y W^\top) - \mathbf{Tr}\left(V^\top W\right) + \epsilon\|W\|_{2,1},$$

with $K^Y = VV^\top$, $W = \Omega V$, $\tilde{K}_t^Y = V^{-1}K_t^Y(V^\top)^{-1}$. Due to the different quadratic terms, this problem cannot be summed up as a Multi-Task Lasso like before. However, it may still be solved, *e.g.* by proximal gradient descent. Indeed, the gradient of the smooth term (*i.e.* all but the $\ell_{2,1}$ mixed norm) reads

$$W + \frac{1}{\Lambda n}\sum_{t=1}^{T}K_t^X W\tilde{K}_t^Y - V, \tag{15}$$

while the proximal operator of the $\ell_{2,1}$ mixed norm is

$$\mathrm{prox}_{\epsilon\,\|\cdot\|_{2,1}}(W) = \begin{pmatrix} | \\ \mathrm{prox}_{\epsilon\,\|\cdot\|_2}(W_{i:}) \\ | \end{pmatrix} = \begin{pmatrix} | \\ \left(1 - \frac{\epsilon}{\|W_{i:}\|_2}\right)_+ W_{i:} \\ | \end{pmatrix} = \begin{pmatrix} | \\ \mathrm{BST}(W_{i:},\epsilon) \\ | \end{pmatrix}.$$

Hence, even in the more involved case of an OVK satisfying only Assumption 4, we have designed an efficient algorithm to compute the solutions to the dual problem.

### A.5.2. $\kappa$-HUBER – FROM PROBLEM $(P2)$ TO $(D2)$

Basic manipulations give the Fenchel-Legendre transforms of the Huber loss:

$$\left(y \mapsto \ell_{H,\kappa}(y - y_i)\right)^\star(\alpha) = \left(\kappa\|\cdot\|_{\mathcal{Y}} \,\square\, \frac{1}{2}\|\cdot\|_{\mathcal{Y}}^2\right)^\star(\alpha) + \langle\alpha, y_i\rangle_{\mathcal{Y}},$$

$$= (\kappa\|\cdot\|_{\mathcal{Y}})^\star(\alpha) + \left(\frac{1}{2}\|\cdot\|_{\mathcal{Y}}^2\right)^\star(\alpha) + \langle\alpha, y_i\rangle_{\mathcal{Y}},$$

$$= \chi_{\mathcal{B}_\kappa}(\alpha) + \frac{1}{2}\|\alpha\|_{\mathcal{Y}}^2 + \langle\alpha, y_i\rangle_{\mathcal{Y}}.$$

Following the same lines as for as for the $\epsilon$-Ridge regression, the dual problem writes

$$\min_{(\alpha_i)_{i=1}^n \in \mathcal{Y}^n}\ \frac{1}{2}\sum_{i,j=1}^{n}\left\langle\alpha_i, \left(\delta_{ij}\mathbf{I}_{\mathcal{Y}} + \frac{1}{\Lambda n}\mathcal{K}(x_i,x_j)\right)\alpha_j\right\rangle_{\mathcal{Y}} - \sum_{i=1}^{n}\langle\alpha_i, y_i\rangle_{\mathcal{Y}} + \sum_{i=1}^{n}\chi_\kappa(\|\alpha_i\|_{\mathcal{Y}}),$$

or again after the reparametrization in $\Omega$

$$\min_{\Omega \in \mathbb{R}^{n \times n}}\ \frac{1}{2}\mathbf{Tr}\left(\tilde{K}\Omega K^Y\Omega^\top\right) - \mathbf{Tr}\left(K^Y\Omega\right)$$

$$\text{s.t.}\quad \sqrt{[\Omega K^Y\Omega^\top]_{ii}} \leqslant \kappa \qquad \forall i \leqslant n$$

The same change of variable permits to conclude. $\qquad\square$

When $\mathcal{K}$ is not identity decomposable, but only satisfies Assumption 4, the problem rewrites

$$\min_{W \in \mathbb{R}^{n \times n}}\ \frac{1}{2}\mathbf{Tr}(WW^\top) + \frac{1}{2\Lambda n}\sum_{t=1}^{T}\mathbf{Tr}(K_t^X W\tilde{K}_t^Y W^\top) - \mathbf{Tr}\left(V^\top W\right),$$

$$\text{s.t.}\quad \|W_{i:}\|_2 \leqslant \kappa \quad \forall i \leqslant n,$$

Again, the gradient term is given by Equation (15), while the projection is similar to the identity decomposable case. The only change thus occurs in the gradient step of Algorithm 1, with a replacement by the above formula.

Notice that if $\kappa$ tends to infinity, the problem is unconstrained, and one also recovers the standard Ridge regression solution.

### A.5.3. $\epsilon$-SVR – FROM PROBLEM $(P3)$ TO $(D3)$

The proof is similar to the above derivations except that the term $\sum_i \|\alpha_i\|_{\mathcal{Y}}^2$ does not appear in the dual, hence the change of matrix $\tilde{K}$. Instead, the dual problem features both the $\ell_{2,1}$ penalization and the $\ell_{2,\infty}$ constraint. $\qquad\square$

### A.6. Proof of Theorem 7

The proof is similar to Appendix A.5.2, with the finite representation coming from Theorem 5.

### A.7. Proof of Theorem 10

In this section, we detail the derivation of constants in Figure 1.

#### A.7.1. $\epsilon$-SVR

Using that the null function is part of the vv-RKHS, it holds

$$\frac{\Lambda}{2}\|h_{A(\mathcal{S})}\|_{\mathcal{H}_\mathcal{K}}^2 \leqslant \hat{\mathcal{R}}_n(h_{A(\mathcal{S})}) \leqslant \hat{\mathcal{R}}_n(0_{\mathcal{H}_\mathcal{K}}) \leqslant M_\mathcal{Y} - \epsilon, \qquad \text{or again} \qquad \|h_{A(\mathcal{S})}\|_{\mathcal{H}_\mathcal{K}} \leqslant \sqrt{\frac{2}{\Lambda}(M_\mathcal{Y} - \epsilon)}.$$

Furthermore, the reproducing property and Assumption 8 give that for any $x \in \mathcal{X}$ and any $h \in \mathcal{H}_\mathcal{K}$ it holds

$$\|h(x)\|^2 = \left\langle \mathcal{K}(\cdot, x)\mathcal{K}(\cdot, x)^\#h, h\right\rangle_{\mathcal{H}_\mathcal{K}} \leqslant \left\|\mathcal{K}(\cdot, x)\mathcal{K}(\cdot, x)^\#\right\|_{\text{op}} \|h\|_{\mathcal{H}_\mathcal{K}}^2 \leqslant \|\mathcal{K}(x,x)\|_{\text{op}} \|h\|_{\mathcal{H}_\mathcal{K}}^2 \leqslant \gamma^2 \|h\|_{\mathcal{H}_\mathcal{K}}^2.$$

Therefore, one gets that for any realization $(x, y) \in \mathcal{X} \times \mathcal{Y}$ of $(X, Y)$ it holds

$$\ell(h_{A(\mathcal{S})}(x), y) = \max(\|y - h_{A(\mathcal{S})}(x)\|_\mathcal{Y} - \epsilon, 0) \leqslant M_y - \epsilon + \|h_{A(\mathcal{S})}(x)\|_\mathcal{Y} \leqslant \sqrt{M_\mathcal{Y} - \epsilon}\left(\gamma\sqrt{\frac{2}{\Lambda}} + \sqrt{M_\mathcal{Y} - \epsilon}\right).$$

This gives $M$. As for $C$, one has

$$\ell(h_{A(\mathcal{S})}(x), y) - \ell(h_{A(\mathcal{S}^{\backslash i})}(x), y) = \max(\|y - h_{A(\mathcal{S})}(x)\|_\mathcal{Y} - \epsilon, 0) - \max(\|y - h_{A(\mathcal{S}^{\backslash i})}(x)\|_\mathcal{Y} - \epsilon, 0).$$

If both norms are smaller than $\epsilon$, then any value of $C$ fits. If both norms are greater than $\epsilon$, the difference reads

$$\|y - h_{A(\mathcal{S})}(x)\|_\mathcal{Y} - \|y - h_{A(\mathcal{S}^{\backslash i})}(x)\|_\mathcal{Y} \leqslant \|h_{A(\mathcal{S})}(x) - h_{A(\mathcal{S}^{\backslash i})}(x)\|_\mathcal{Y}.$$

If only one norm is greater than $\epsilon$ (we write it only for $h_{A(\mathcal{S})}$ as it is symmetrical), the difference may be rewritten

$$\|y - h_{A(\mathcal{S})}(x)\|_\mathcal{Y} - \epsilon \leqslant \|y - h_{A(\mathcal{S})}(x)\|_\mathcal{Y} - \|y - h_{A(\mathcal{S}^{\backslash i})}(x)\|_\mathcal{Y} \leqslant \|h_{A(\mathcal{S})}(x) - h_{A(\mathcal{S}^{\backslash i})}(x)\|_\mathcal{Y}.$$

Hence we get $C = 1$.

#### A.7.2. $\epsilon$-RIDGE

Using the same reasoning as for the $\epsilon$-SVR, one has

$$\|h_{A(\mathcal{S})}\|_{\mathcal{H}_\mathcal{K}} \leqslant \sqrt{\frac{2}{\Lambda}(M_\mathcal{Y} - \epsilon)} \qquad \text{and} \qquad \|h_{A(\mathcal{S}^{\backslash i})}\|_{\mathcal{H}_\mathcal{K}} \leqslant \sqrt{\frac{2}{\Lambda}(M_\mathcal{Y} - \epsilon)}. \tag{16}$$

Therefore, for any realization $(x, y) \in \mathcal{X} \times \mathcal{Y}$ of $(X, Y)$ it holds

$$\ell(h_{A(\mathcal{S})}(x), y) = \max(\|y - h_{A(\mathcal{S})}(x)\| - \epsilon, 0)^2 \leqslant (\|y\|_\mathcal{Y} - \epsilon + \|h_{A(\mathcal{S})}(x)\|_\mathcal{Y})^2 \leqslant (M_\mathcal{Y} - \epsilon)^2\left(1 + \frac{2\sqrt{2}\gamma}{\sqrt{\Lambda}} + \frac{2\gamma^2}{\Lambda}\right).$$

As for $C$, one has

$$\ell(h_{A(\mathcal{S})}(x), y) - \ell(h_{A(\mathcal{S}^{\backslash i})}(x), y) = \max(\|y - h_{A(\mathcal{S})}(x)\|_\mathcal{Y} - \epsilon, 0)^2 - \max(\|y - h_{A(\mathcal{S}^{\backslash i})}(x)\|_\mathcal{Y} - \epsilon, 0)^2.$$

If both norms are smaller than $\epsilon$, any $C$ fits. If both are larger than $\epsilon$, using Equation (16) the difference becomes

$$\left(\|y - h_{A(\mathcal{S})}(x)\|_\mathcal{Y} + \|y - h_{A(\mathcal{S}^{\backslash i})}(x)\|_\mathcal{Y} - 2\epsilon\right)\left(\|y - h_{A(\mathcal{S})}(x)\|_\mathcal{Y} - \|y - h_{A(\mathcal{S}^{\backslash i})}(x)\|_\mathcal{Y}\right),$$

$$\leqslant 2(M_\mathcal{Y} - \epsilon)\left(1 + \frac{\gamma\sqrt{2}}{\sqrt{\Lambda}}\right)\|h_{A(\mathcal{S})}(x) - h_{A(\mathcal{S}^{\backslash i})}(x)\|_\mathcal{Y}.$$

If only one norm is greater than $\epsilon$ (again, the analysis is symmetrical), the difference may be rewritten

$$\left(\|y - h_{A(\mathcal{S})}(x)\|_{\mathcal{Y}} - \epsilon\right)^2 \leqslant \left(\|y - h_{A(\mathcal{S})}(x)\|_{\mathcal{Y}} - \|y - h_{A(\mathcal{S}^{\backslash i})}(x)\|_{\mathcal{Y}}\right)^2 \leqslant \|h_{A(\mathcal{S})}(x) - h_{A(\mathcal{S}^{\backslash i})}(x)\|_{\mathcal{Y}}^2,$$

$$\leqslant \left(\|h_{A(\mathcal{S})}(x)\|_{\mathcal{Y}} + \|h_{A(\mathcal{S}^{\backslash i})}(x)\|_{\mathcal{Y}}\right) \|h_{A(\mathcal{S})}(x) - h_{A(\mathcal{S}^{\backslash i})}(x)\|_{\mathcal{Y}},$$

$$\leqslant 2(M_{\mathcal{Y}} - \epsilon)\frac{\gamma\sqrt{2}}{\sqrt{\Lambda}}\|h_{A(\mathcal{S})}(x) - h_{A(\mathcal{S}^{\backslash i})}(x)\|_{\mathcal{Y}}.$$

In every case $C = 2(M_{\mathcal{Y}} - \epsilon)\left(1 + \gamma\sqrt{2}/\sqrt{\Lambda}\right)$ works, hence the conclusion.

### A.7.3. $\kappa$-HUBER

Using the same techniques, one gets

$$\|h_{A(\mathcal{S})}\|_{\mathcal{H}_{\mathcal{K}}} \leqslant \sqrt{\frac{2\kappa}{\Lambda}\left(M_{\mathcal{Y}} - \frac{\kappa}{2}\right)} \qquad \text{and} \qquad \|h_{A(\mathcal{S}^{\backslash i})}\|_{\mathcal{H}_{\mathcal{K}}} \leqslant \sqrt{\frac{2\kappa}{\Lambda}\left(M_{\mathcal{Y}} - \frac{\kappa}{2}\right)},$$

and for any realization $(x, y) \in \mathcal{X} \times \mathcal{Y}$ of $(X, Y)$

$$\ell(h_{A(\mathcal{S})}(x), y) \leqslant \kappa\sqrt{M_{\mathcal{Y}} - \frac{\kappa}{2}}\left(\frac{\gamma\sqrt{2\kappa}}{\sqrt{\Lambda}} + \sqrt{M_{\mathcal{Y}} - \frac{\kappa}{2}}\right).$$

If both norms are greater than $\kappa$, the difference $\ell(h_{A(\mathcal{S})}(x), y) - \ell(h_{A(\mathcal{S}^{\backslash i})}(x), y)$ writes

$$\kappa\left(\|h_{A(\mathcal{S})}(x) - y\|_{\mathcal{Y}} - \frac{\kappa}{2}\right) - \kappa\left(\|h_{A(\mathcal{S}^{\backslash i})}(x) - y\|_{\mathcal{Y}} - \frac{\kappa}{2}\right) \leqslant \kappa\|h_{A(\mathcal{S})}(x) - h_{A(\mathcal{S}^{\backslash i})}(x)\|_{\mathcal{Y}}.$$

If only one norm is greater than $\kappa$, one may upperbound the difference using the previous writing

$$\kappa\left(\|h_{A(\mathcal{S})}(x) - y\|_{\mathcal{Y}} - \frac{\kappa}{2}\right) - \frac{1}{2}\|h_{A(\mathcal{S}^{\backslash i})}(x) - y\|_{\mathcal{Y}}^2 \leqslant \kappa\left(\|h_{A(\mathcal{S})}(x) - y\|_{\mathcal{Y}} - \frac{\kappa}{2}\right) - \kappa\left(\|h_{A(\mathcal{S}^{\backslash i})}(x) - y\|_{\mathcal{Y}} - \frac{\kappa}{2}\right).$$

If both are smaller than $\kappa$, the difference becomes

$$\frac{1}{2}\|h_{A(\mathcal{S})}(x) - y\|_{\mathcal{Y}}^2 - \frac{1}{2}\|h_{A(\mathcal{S}^{\backslash i})}(x) - y\|_{\mathcal{Y}}^2,$$

$$= \frac{1}{2}\left(\|h_{A(\mathcal{S})}(x) - y\|_{\mathcal{Y}} + \|h_{A(\mathcal{S}^{\backslash i})}(x) - y\|_{\mathcal{Y}}\right)\left(\|h_{A(\mathcal{S})}(x) - y\|_{\mathcal{Y}} - \|h_{A(\mathcal{S}^{\backslash i})}(x) - y\|_{\mathcal{Y}}\right),$$

$$\leqslant \kappa\|h_{A(\mathcal{S})}(x) - h_{A(\mathcal{S}^{\backslash i})}(x)\|_{\mathcal{Y}},$$

so that $C = \kappa$.

## A.8. Further Admissible Kernels for Assumption 3

In the continuation of Remark 1, we now exhibit several types of OVK that satisfy Assumption 3.

**Proposition 2.** *The following Operator-Valued Kernels satisfy Assumption 3:*

*(i)* $\forall s, t \in \mathcal{X}^2, \ \mathcal{K}(s, t) = \sum_i k_i(s, t) \, y_i \otimes y_i,$      *with $k_i$ positive semi-definite (p.s.d.) scalar kernels for all $i \leqslant n$.*

*(ii)* $\forall s, t \in \mathcal{X}^2, \ \mathcal{K}(s, t) = \sum_i \mu_i \, k(s, t) \, y_i \otimes y_i,$      *with $k$ a p.s.d. scalar kernel and $\mu_i \geqslant 0$ for all $i \leqslant n$.*

*(iii)* $\forall s, t \in \mathcal{X}^2, \ \mathcal{K}(s, t) = \sum_i k(s, x_i)k(t, x_i) \, y_i \otimes y_i,$

*(iv)* $\forall s, t \in \mathcal{X}^2, \ \mathcal{K}(s, t) = \sum_{i,j} k_{ij}(s, t) \, (y_i + y_j) \otimes (y_i + y_j),$      *with $k_{ij}$ p.s.d. scalar kernels for all $i, j \leqslant n$.*

*(v)* $\forall s, t \in \mathcal{X}^2, \ \mathcal{K}(s, t) = \sum_{i,j} \mu_{ij} \, k(s, t) \, (y_i + y_j) \otimes (y_i + y_j),$      *with $k$ a p.s.d. scalar kernel and $\mu_{ij} \geqslant 0$.*

*(vi)* $\forall s, t \in \mathcal{X}^2, \ \mathcal{K}(s, t) = \sum_{i,j} k(s, x_i, x_j)k(t, x_i, x_j) \, (y_i + y_j) \otimes (y_i + y_j).$

*Proof.*

(i) For all $(s_k, z_k)_{k \leqslant n} \in (\mathcal{X} \times \mathcal{Y})^n$, it holds: $\sum_{k,l} \langle z_k, \mathcal{K}(s_k, s_l) z_k \rangle_{\mathcal{Y}} = \sum_i \sum_{k,l} k_i(s, t) \langle z_k, y_i \rangle_{\mathcal{Y}} \langle z_l, y_i \rangle_{\mathcal{Y}}$, which is positive by the positiveness of the scalar kernels $k_i$'s. Notice that $(ii)$ and $(iii)$ are then particular cases of $(i)$.

(ii) is an application of $(i)$, as a kernel remains p.s.d. through positive multiplication. Observe that this kernel is separable.

(iii) is also a direct application of $(i)$, kernel $k' : s, t \mapsto k(s, x_i) k(t, x_i)$ being indeed p.s.d. for all function $k$ and point $x_i$.

(iv) is proved similarly to $(i)$. The arguments used for $(ii)$ and $(iii)$ also makes $(v)$ and $(vi)$ direct applications of $(iv)$.

Finally, notice that for $(iv)$, $(v)$ and $(vi)$, any linear combination $(\nu_i y_i + \nu_j y_j) \otimes (\nu_i y_i + \nu_j y_j)$, with $0 \leqslant \nu_i \leqslant 1$ for all $i \leqslant n$, could have been used instead of $(y_i + y_j) \otimes (y_i + y_j)$. $\qquad \square$

## B. Loss Functions Illustrations

In this section, we provide illustrations of the loss functions we used to promote sparsity and robustness. This includes $\epsilon$-insensitive losses (Definitions 3 and 4, Figures 9 and 10) and the $\kappa$-Huber loss (Definition 5, Figure 11). First introduced for real outputs, their formulations as infimal convolutions allows for a generalization to any Hilbert space, either of finite dimension (as in Sangnier et al. (2017)) or not, which is the general case addressed in the present paper. The $\epsilon$-insensitive loss functions promote sparsity, as reflected in the corresponding dual problems (see Theorem 6, Problems $(D1)$ and $(D3)$ therein) and the empirical results (Figures 12 and 13). On the other hand, losses whose slopes asymptotically behave as $|| \cdot ||_{\mathcal{Y}}$ instead of $|| \cdot ||_{\mathcal{Y}}^2$ (such as the $\kappa$-Huber or the $\epsilon$-SVR loss) encourage robustness through a resistance to outliers. Indeed, under such a setting, residuals of high norm contribute less to the gradient and have a minor influence on the model output.
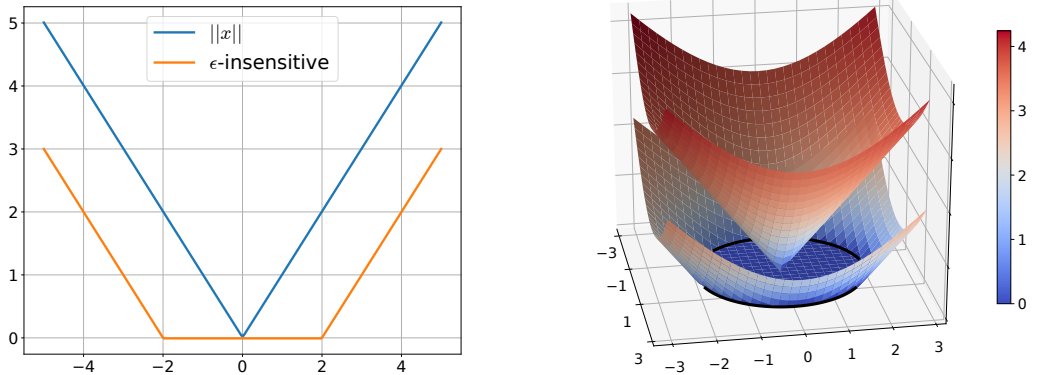


*Figure 9.* Standard and $\epsilon$-insensitive versions of the SVR loss in 1 and 2 dimensions ($\epsilon = 2$).
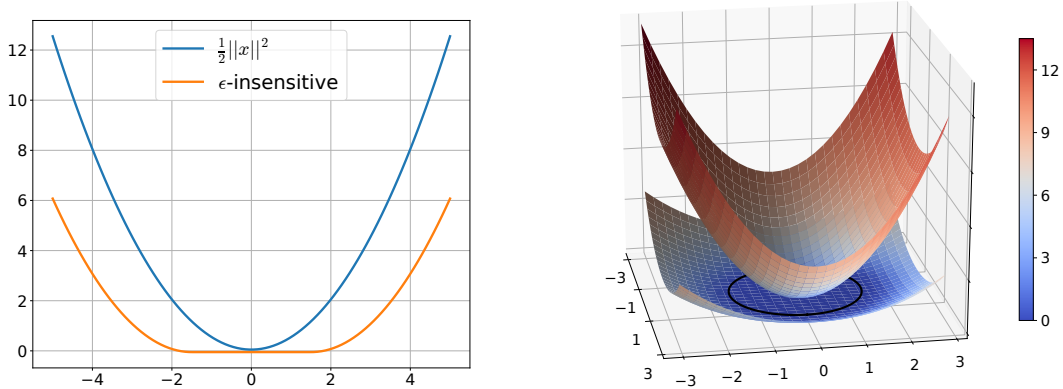


*Figure 10.* Standard and $\epsilon$-insensitive versions of the square loss in 1 and 2 dimensions ($\epsilon = 1.5$).
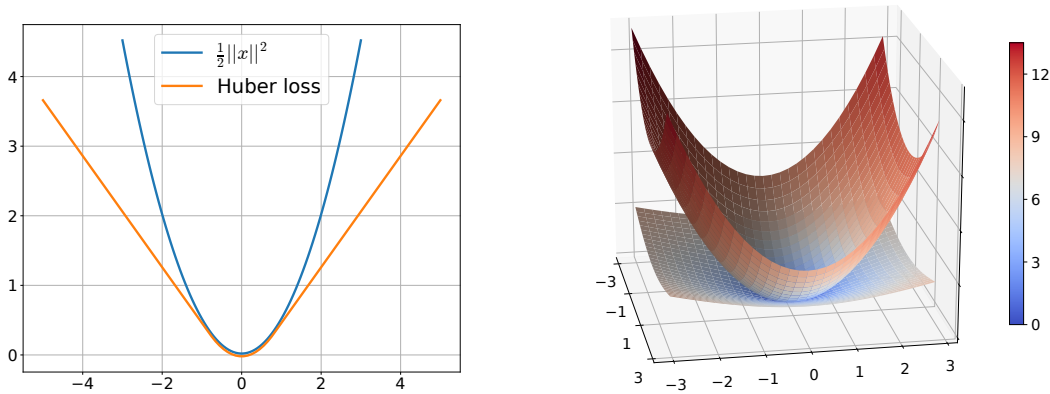
*Figure 11.* Standard square loss and Huber loss in 1 and 2 dimensions ($\kappa = 0.8$).

# C. Numerical Experiments and Code

## C.1. Provided Code

The Python code used to generate the plots and tables of the article is provided. The `README` file in the code folder contains instructions for quickly reproducing (part of) the plots. All implemented methods may be run on other datasets/problems.

## C.2. Detailed Protocols

### C.2.1. STRUCTURED PREDICTION

**YEAST Dataset Description.** YEAST[1] is a publicly available multi-label classification dataset used as a benchmark in several structured prediction articles. We compared our approach, with the same train/test decomposition, to those presented in Elisseeff and Weston (2002), Finley and Joachims (2008) and Belanger and McCallum (2016). The size of the training set is 1500, the test set is of size 917. The problem consists in predicting the functional classes of a gene. The inputs are micro-array expression data (representing the genes) of dimension $p = 103$. The outputs are multi-label vectors of size $d = 14$ representing the possible functional classes of the genes. The average number of labels is 4.2. These 14 functional classes correspond to the first level of a tree that structures a much bigger set of possible functional classes.

**Experimental protocol: Comparison with other methods.** In Figure 6, we reported the Hamming error on the test set obtained by each method. The results obtained by SSVM and SPENS are extracted from Finley and Joachims (2008) and Belanger and McCallum (2016). For our approach and its three variants ($\epsilon$-KRR, $\kappa$-Huber, $\epsilon$-SVR), each hyper-parameter ($\Lambda$, $\epsilon$, or $\kappa$) has been selected by estimating the Mean Squared Error (MSE) through a 5-fold cross-validation computed on the training set. We used an input Gaussian kernel with a fixed bandwidth equal to 1.

**Experimental protocol: Cross-Effect of $\epsilon$ and $\Lambda$ on sparsity and MSE.** In order to measure the effect of the different hyperparameters and study their interrelations, we have computed the 5-fold cross-validation MSE and sparsity/saturation for several values of $\Lambda$ and $\epsilon/\kappa$. The input kernel is still Gaussian with bandwidth 1. The results are plotted in Figures 3 and 4 for the $\epsilon$-KRR, and in Figures 13 and 14 for the $\epsilon$-SVR and $\kappa$-Huber. In Figure 4, we have measured sparsity through the number of training data which are discarded, *i.e.* not used in the finite representation of the $\epsilon$-KRR model. The $\kappa$-Huber saturation is assessed in a similar fashion: it corresponds to the number of training data whose associated coefficient saturates the norm constraint (see Theorem 6, Problem ($D2$) therein). Simplified versions of these graphs may be quickly reproduced using the code attached (see `README` file).

**Metabolite identification dataset description.** We next tested our method on a harder problem: that of metabolite identification (Brouard et al., 2016a). The goal is to predict a metabolite (small molecule) thanks to its mass spectrum. The difficulty comes from the reduced size of the training set ($n = 6974$) compared to the high dimension of the outputs ($d = 7593$). Input Output Kernel Regression (IOKR, see Brouard et al. (2016a;b)) with a Tanimoto-Gaussian kernel is state of the art on this problem.

---

[1] https://www.csie.ntu.edu.tw/ cjlin/libsvmtools/datasets/

**Experimental protocol.** We investigate the advantages of substituting the Ridge Regression for the $\epsilon$-KRR, $\kappa$-Huber, and $\epsilon$-SVR. Outputs are embedded in an infinite dimensional space through the use of the Tanimoto-Gaussian kernel (with bandwidth $\gamma = 0.72$). We compare the different algorithms' performances on a set of 6974 mass spectra through the top-$k$ accuracies for $k \in \{1, 10, 20\}$. We give the average 5-fold top-$k$ accuracies (Table 1). The 5 folds have been chosen such that a metabolite does not appear in two different folds (zero-shot learning setting).

### C.2.2. STRUCTURED REPRESENTATION LEARNING

**Dataset Description.** Robust structured representation learning was tested on a drug dataset, introduced in Su et al. (2010), and extracted from the NCI-Cancer database. This dataset features a set of molecules that are represented through a Gram matrix of size $2303 \times 2303$ obtained with a Tanimoto kernel. Tanimoto kernels (see Ralaivola et al. (2005) for details) are a common way to compare labeled graphs by means of a bag-of-sequences approach.

**Experimental protocol: Robust KAE.** We computed the mean 5-fold cross-validation Mean Squared Error. The first layer uses a linear kernel. But since inputs (and outputs) are kernelized – only the $2303 \times 2303$ Gram matrix is provided for learning –, the first layer may also be seen as a function from the associated Tanimoto-RKHS, applied to the molecules. The second layer uses a Gaussian kernel. The regularization parameters for the two layers have been fixed to $\Lambda = 1e - 6$, and the inner dimension has been set to $p = 200$. In Figure 7 is plotted the MSE and the sparsity (discarded training data) for several values of $\epsilon$ in order to assess the effect of the regularization. We used an existing source code from Laforgue et al. (2019)[2], that has been adapted to our needs. The IOKR resolution part, materialized by the `compute_N_L` function therein, has been replaced by the `compute_Omega` function of the `IOKR_plus` class in the attached code.

### C.2.3. ROBUST FUNCTION-TO-FUNCTION REGRESSION

**Dataset Description.** The task at hand consists in predicting lip acceleration from electromyography (EMG) signals of the corresponding muscle (Ramsay and Silverman, 2007). The dataset[3] includes 32 samples of time series obtained by recording a subject saying "say bob again", that are noted $(x_i, y_i)_{i=1}^{32}$. Each time series is of length 64. To assess the performance of our method in presence of outliers, we created 4 outliers by picking randomly some $(x_i)_{i=1}^{4}$ and adding to the dataset the samples $(x_i, -1.2 * y_i)_{i=1}^{4}$.

**Experimental protocol.** As the number of samples is small, one can use the Leave One Out (LOO) generalization error as a measure of the model performance. We first used it with plain Ridge Regression (Kadri et al., 2016) to select the best hyperparameter $\Lambda$. Then, we tested robustness by computing the LOO generalization error of a model output by solving Problem (9) for various $\kappa$ (see Figure 8, that may also be reproduced from the attached code). For the $\{\psi_j\}_{j=1}^{m}$ we used the sine and cosine basis of $L^2([0,1])$, i.e. $\forall l \leqslant \frac{m}{2}$ and $\theta \in [0,1]$, $\psi_{2l}(\theta) = \sqrt{2}\cos(2\pi l\theta)$ and $\psi_{2l+1}(\theta) = \sqrt{2}\sin(2\pi l\theta)$. The number of basis function was set to $m = 16$, so that we get the first 8 cosines and sines of the basis. The chosen associated eigenvalues are $\lambda_{2l} = \lambda_{2l+1} = \frac{1}{(1+j)^2}$. We used as an input kernel the integral Laplacian $k_{\mathcal{X}}(x_1, x_2) = \int_0^1 \exp\left(-7|x_1(\theta) - x_2(\theta)|\right)\mathrm{d}\theta$.

### C.3. Additional Figures

We now provide analogues to Figures 3 and 4 for the $\epsilon$-SVR and $\kappa$-Huber. The $\epsilon$-Ridge graphs are first recalled. Notice that simplified versions of these plots may be easily reproduced from the attached code.

The $\epsilon$-KRR (Figure 12) appears as a natural regularized version of the plain KRR. For small values of $\Lambda$, the regularization effect of the $\epsilon$ induces a smaller MSE. This phenomenon is achieved for a wide range of $\Lambda$ and $\epsilon$, and coincides with an important sparsity. The counterpart is that no value of $\epsilon$ clearly allows to outperform the standard KRR for its optimal $\Lambda$. The $\epsilon$-KRR may rather be used as an implicit regularization preventing from a cross-validation on $\Lambda$.

The $\epsilon$-SVR (Figure 13) shares analogous characteristics for the small $\Lambda$ regime. However, it further produces predictors with smaller MSE than the best KRR one. This furthermore coincides with a peak in the sparsity.

The $\kappa$-Huber (Figure 14) has a quite different behavior. When $\Lambda$ tends to 0, the constraint (see Problem ($D2$)) is vacuous for all $\kappa$, and one asymptotically recovers the standard KRR. The optimal $\Lambda$ now changes with $\kappa$, and better performances than the KRR for the best $\Lambda$ are regularly attained.

---

[2]github.com/plaforgue/kae
[3]http://www.stats.ox.ac.uk/ silverma/fdacasebook/lipemg.html

*Figure 12.* MSE and Sparsity w.r.t. $\Lambda$ for different $\epsilon$ for the $\epsilon$-KRR on the YEAST dataset.



*Figure 13.* MSE and Sparsity w.r.t. $\Lambda$ for different $\epsilon$ for the $\epsilon$-SVR on the YEAST dataset.



*Figure 14.* MSE and Saturation w.r.t. $\Lambda$ for different $\kappa$ for the $\kappa$-Huber on the YEAST dataset.

REFERENCES

Belanger, D. and McCallum, A. (2016). Structured prediction energy networks. In *International Conference on Machine Learning*, pages 983–992.

Brouard, C., Shen, H., Dührkop, K., d'Alché Buc, F., Böcker, S., and Rousu, J. (2016a). Fast metabolite identification with input output kernel regression. *Bioinformatics*, 32(12):i28–i36.

Brouard, C., Szafranski, M., and D'Alché-Buc, F. (2016b). Input output kernel regression: supervised and semi-supervised structured output prediction with operator-valued kernels. *The Journal of Machine Learning Research*, 17(1):6105–6152.

Elisseeff, A. and Weston, J. (2002). A kernel method for multi-labelled classification. In *Advances in neural information processing systems*, pages 681–687.

Finley, T. and Joachims, T. (2008). Training structural svms when exact inference is intractable. In *Proceedings of the 25th international conference on Machine learning*, pages 304–311.

Kadri, H., Duflos, E., Preux, P., Canu, S., Rakotomamonjy, A., and Audiffren, J. (2016). Operator-valued kernels for learning from functional response data. *Journal of Machine Learning Research*, 17(20):1–54.

Laforgue, P., Clémençon, S., and d'Alché-Buc, F. (2019). Autoencoding any data through kernel autoencoders. In *Artificial Intelligence and Statistics*, pages 1061–1069.

Ralaivola, L., Swamidass, S. J., Saigo, H., and Baldi, P. (2005). Graph kernels for chemical informatics. *Neural networks*, 18(8):1093–1110.

Ramsay, J. O. and Silverman, B. W. (2007). *Applied functional data analysis: methods and case studies*. Springer.

Sangnier, M., Fercoq, O., and d'Alché-Buc, F. (2017). Data sparse nonparametric regression with $\epsilon$-insensitive losses. In *Asian Conference on Machine Learning*, pages 192–207.

Su, H., Heinonen, M., and Rousu, J. (2010). Structured output prediction of anti-cancer drug activity. In *IAPR International Conference on Pattern Recognition in Bioinformatics*, pages 38–49. Springer.